

Deep QWOP Learning

Hung-Wei Wu

Submitted under the supervision of Maria Gini and James Parker to the University Honors Program at the University of Minnesota-Twin Cities in partial fulfillment of the requirements for the degree of Bachelor of Sciences *cum laude* in Computer Science.

12/2/2017

1 Abstract

We apply a deep learning model to the QWOP flash game, which requires control of a ragdoll athlete using only the keys “Q”, “W”, “O”, and “P”. The model is a convolutional neural network trained with Q-learning. By training the model with only raw pixel input, we show that our model is capable of successfully learning a control policy associated with playing QWOP. This model was successfully applied to a non-deterministic control environment in the form of a ragdoll physics flash game.

2 Introduction

2.1 QWOP

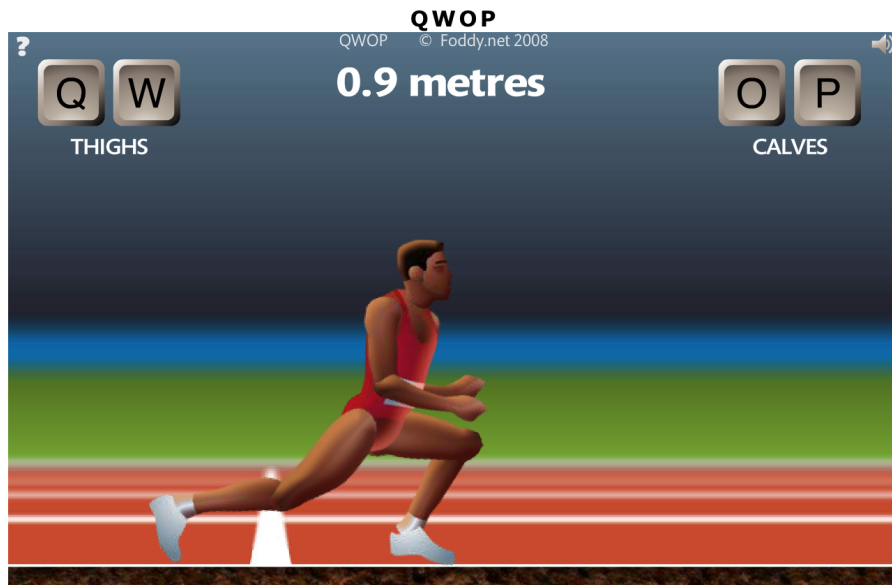


Figure 1. QWOP game play.

QWOP is a free-to-play flash game created by Bennet Foddy infamous for being ridiculously frustrating to play [10]. In QWOP, the user controls a ragdoll sprinter using the four keys: “Q”, “W”, “O”, and “P”. Each key controls the left thigh, left calf, right thigh, and right calf respectively. With the right inputs and timing, this can be used to simulate real-world

human-like running. However, this is not how we as humans, are used to running. Our motor skills usually don't involve thinking about how specific muscles have to move in order to move forward and maintain balance. This means that in the context of QWOP, the player's collective knowledge on balance and movement is essentially useless [9]. The goal is for the user to attempt to move the ragdoll figure 100 meters without falling over. The game is reset when any section of the upper torso touches the ground.

The game implements a ragdoll physics environment where complicated interactions such as gravity and momentum are greatly simplified as a tradeoff for low CPU utilization when rendering. In particular, this means that any body parts that are not being directly simulated is latent, meaning it just falls in the direction that it is already traveling. If the runner gets slightly out of balance and without the player's intervention, it will fall. The articulated figure has little to zero joint stiffness, often leading to it collapsing into comically improbable or compromising positions. The game is notoriously difficult and achieving any sort of forward movement is considered a significant achievement.

2.2 Deep Q Learning

DeepMind published a paper in 2013 "Playing Atari with Deep Reinforcement Learning" describing a deep reinforcement learning system that combines neural networks with reinforcement learning to master a diverse range of Atari 2600 games using only the raw pixels and score as inputs [6]. Until this point, it has only been possible to create individual algorithms capable of mastering a single specific domain [13]. Deep Q Learning represents the first demonstration of a general-purpose agent that is able to continually adapt its behavior without human intervention [5]. However, it has only been applied to deterministic tasks, where a given action produces a given result that can be inferred from the environment [15]. The task of

playing QWOP poses a different type of problem. It is significantly more difficult due to the ragdoll physics environment. Each key press is not guaranteed to have the same results or effects on the simulation. Miniscule differences in the runner's position and momentum can often have unforeseen impacts.

3 Related Work

3.1 DeepMind Atari

Google DeepMind published a paper in 2013 describing the first deep learning model to successfully learn control policies directly from sensory input using reinforcement learning [6]. The input is raw pixels and the output is a value function estimating future rewards. Their method was able to learn to play seven Atari 2600 games and even surpass a human expert on three of the games. These games include Pong, Breakout, Space Invaders, Seaquest, and Beam Rider. Their model is a convolutional neural network trained with a variant of Q-learning, using stochastic gradient descent to update the weights. They also implemented an experience replay mechanism which randomly samples previous actions and state transitions to smooth out the training distribution over past behaviors [3]. Our model is based on this architecture, we will be implementing a convolutional neural network trained with Q-learning.

3.2 OpenAI Gym

OpenAI Gym is a toolkit for developing and comparing reinforcement learning algorithms and techniques [7]. This platform provides many environments that agents can interact with in a unified way. It provides an interface that allows agents to step the environment by one timestep and return new observations, rewards, and exit statuses.

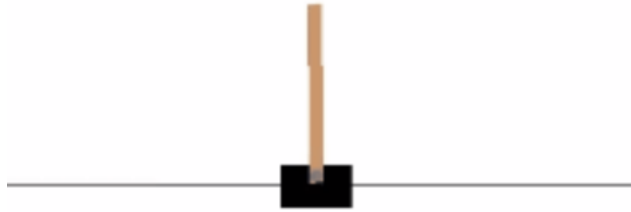


Figure 2. OpenAI CartPole environment.

For example, in the CartPole environment, a pole is attached by an un-actuated joint to a cart, which moves along a frictionless track. This system is controlled by applying a force of +1 or -1 corresponding to left and right movement to the cart. The pendulum initially starts upright and the goal is to prevent it from falling over. A reward of +1 is provided for every timestep that the pole remains upright. The current episode ends when the pole is more than 15 degrees from vertical or when the cart moves more than 2.4 units from the center. CartPole is one of the simplest environments in OpenAI gym. An agent can move the cart by performing a series of actions of 0 or 1 to the cart, pushing it left or right. The QWOP game interface is written to follow a similar environment architecture where an agent has access to methods that allow the it to reset the environment as well as execute actions. An example agent found in the documentation implemented a simple three-layer convolutional neural network and is trained using Q-learning. After around 500 episodes, the agent learned how to maximize the score by keeping the pole upright and the cart in the center of the environment. It is then consistently able to survive all 500 timesteps in each episode.

3.3 Stanford CS229

Gustav Brodman and Ryan Voldstad used reinforcement learning to play QWOP for their CS229 final project [9]. Methods included discretization of state spaces with both regular and

fitted value iteration using a set of reward features. Instead of using raw pixel inputs, other variables were used to better quantify the QWOP runner's state. Distance alone was not enough to determine the state of the runner; therefore, other variables such as number of feet on the ground, left and right knee angles, angle between the left and right legs, and thigh rotational velocities were used to represent the state instead. Through some experimentation, they settled on a feature mapping using the difference between thigh angles, the angles of each knee, the overall "tilt" of the runner, and the runner's horizontal speed. Evaluating their model showed fairly good results. The QWOP sprinter was able to travel around 30000 units (arbitrary distance units). Initially, a shuffling gait was observed; however, after 10 iterations, a gait that resembled bipedal walking was observed.

4 Background

Our QWOP agent implements the Deep Q Learning algorithm using a neural net and reinforcement learning.

4.1 Markov Decision Processes

Markov decision processes provide a mathematical framework for modeling decision-making in situations where outcomes are partly random and partly under the control of a decision maker [11]. At each timestep, a Markov decision process is in some state "s", and the decision maker may choose any action "a" that is available in that particular state. The process responds at the next timestep by randomly moving into a new state "s'", and giving the decision maker a corresponding reward. We are attempting to model QWOP as a Markov decision process even though identical actions in the same state may not have the same results. The momentum of the ragdoll runner is not captured in the raw pixel input.

4.2 Reinforcement Learning

General reinforcement learning is an area of machine learning inspired by behaviorist psychology. It addresses problems concerning how agents should take actions in an environment to maximize some predefined reward. Reinforcement learning differs from standard supervised learning in that sub-optimal actions are not explicitly corrected, nor correct input and output pairs ever presented [13]. It instead focuses on finding a balance between exploration and usage of current knowledge. In general, an agent performs some action “A” that results in a new state “S” and reward “R”, this is then fed back into the agent. Reinforcement learning is relevant to an enormous range of tasks, including robots, game playing, consumer modeling, and healthcare.

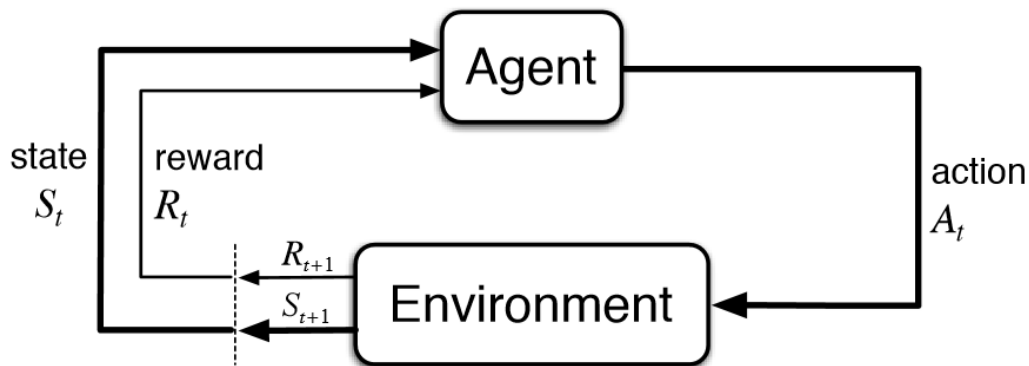


Figure 3. Reinforcement learning architecture.

4.2 Q-Learning

Q-learning is a model-free reinforcement learning technique. Specifically, Q-learning can be used to find an optimal action-selection policy for any given finite Markov decision process. A policy is a rule that the agent follows when selecting actions. In Q-learning, there is an action-value function called the Q-function, which is used to approximate the reward based on a state [2]. It ultimately gives the expected utility of taking a given action in a given state and following the optimal policy thereafter. When such an action-value function is learned, the optimal policy

can be constructed by simply selecting the highest values in each state. One of the strengths of Q-learning is that it is able to compare the expected utilities of the available actions without requiring a model of the environment [4]. Additionally, Q-learning can handle problems with stochastic transitions and rewards, without requiring any adaptations [14]. It has been proven that for any finite Markov decision process, Q-learning eventually finds an optimal policy [1]. We use a convolutional neural network to model the Q-function. The loss function used to train the network is shown below in Figure 4.

$$loss = \left(\underbrace{r + \gamma \max_{a'} \hat{Q}(s, a')}_{\text{Target}} - \underbrace{Q(s, a)}_{\text{Prediction}} \right)^2$$

Reward
Decay Rate

Figure 4. Q value and loss calculation

An agent first carries out an action "a" and observes the reward "r" and the resulting state "s' ". Based on the result, we calculate the maximum target Q-value and then discount it so that the future reward is worth less than the immediate reward.

4.3 Convolutional Neural Networks

A regular neural network receives a single vector as input and transforms it through a series of hidden layers, made of a set of neurons. Each neuron is fully connected to all neurons in the previous layer. Neurons in a single layer function completely independent of each other. The last layer of a network is called the output layer and in classification settings, it represents the class scores.

Convolutional neural networks take advantage of the fact that the input consists of images and thus it constrains the architecture in a more sensible way [3]. In particular, unlike a regular neural network, the layers of a convolutional neural network have neurons arranged in

three dimensions: width, height, and depth. The neurons in a layer will only be connected to a small region of the layer before it, instead of all the neurons in a fully-connected manner. This architecture is visualized below in Figure 5 and Figure 6. We will be using Keras, which is a Python deep learning library [8] to build our convolutional network. The Q-function is modeled using this network.

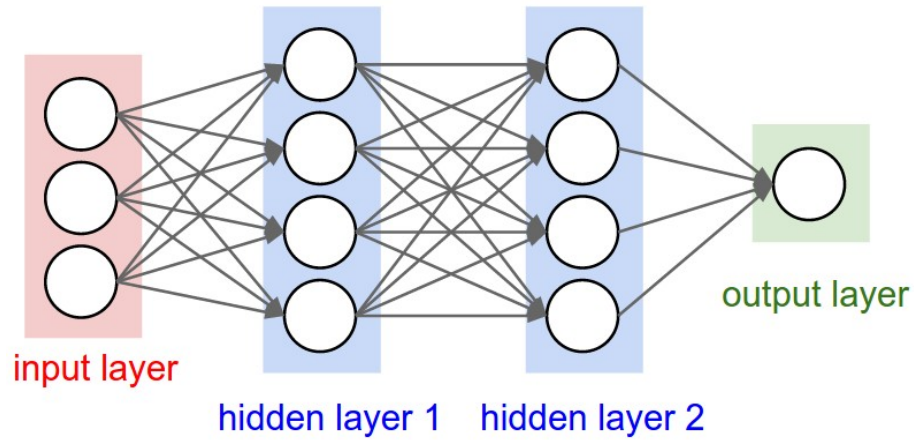


Figure 5. Regular neural network architecture.

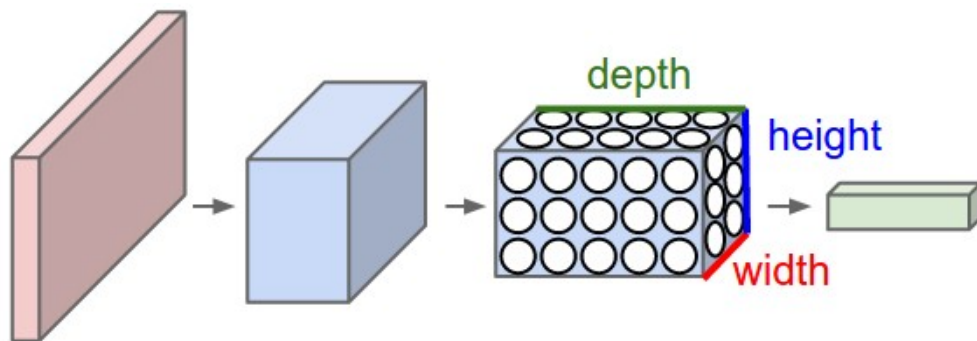


Figure 6. Convolutional neural network architecture.

4.4 Remember and Replay

The most notable features of the Deep Q Learning algorithm are the "remember" and "replay" methods. One of the challenges of Deep Q Learning is that the neural network used in the algorithm tends to forget the previous experiences as it overwrites them with new

experiences [6]. Thus, methods are needed to remember previous actions and rewards and retrain the neural network to retain previous knowledge. To ensure the agent performs well long term, we need to take into account the immediate and future rewards. In order to accomplish this, a discount rate is specified. Thus, the agent will learn to maximize the discounted future reward based on the given state.

4.5 Hyperparameters

There are also some hyperparameters that have to be specified when the model is being trained. They are listed below in Figure 7. The episode parameter specifies how many “games” the agent will play. Each episode has 500 timesteps or actions.

The exploration rate is specified by epsilon. Initially, the neural network is not trained to maximize the Q-function. Thus, the QWOP agent will randomly select possible actions a set percentage of the time. This percentage is specified by the exploration rate. It is better for the agent to try different actions and observe the subsequent rewards and start converging on the optimal action-value function. However, when the agent is not randomly deciding its actions, it will predict the reward value based on the current state and pick the action that will give the highest reward. The exploration rate starts at 1.0 and will gradually decrease over time.

Learning rate in the context of neural networks is a measure of how quickly a network abandons old beliefs for new ones. Neural networks are often trained by gradient descent on the weights. This means that at each iteration we use backpropagation to calculate the derivative of the loss function with respect to each weight and subtract it from that weight. However, in practice, if this is applied, the weights will vary too much and overcorrect and the loss will diverge [3]. Thus, the learning rate is a small value that acts as a multiplier to the derivative of the loss function.

Episodes	The number of games the agents are going to play.
Gamma	The decay rate used to calculate the future discounted reward.
Epsilon	The percentage that the agent will randomly decide its actions.
Epsilon decay	As the network gradually learns patterns, it will explore less and less.
Learning rate	How much the network learns in each iteration.

Figure 7. Deep Q Learning hyperparameters.

4.6 ReLu

The two hidden layers in the neural network used to train the Q-function are composed of rectified linear unit neurons (ReLU). The ReLu is an activation function defined as the positive part of its argument. The function is shown below in Figure 8, where x is the input to a neuron. It was first introduced in 2000 with strong biological motivations and mathematical justifications. It has been used in convolutional networks more effectively than the widely used logistic sigmoid function. ReLu neurons are faster to compute since they do not require any normalization. They also do not require any exponential computation such as those required in sigmoid or tanh activation functions [12]. However, it is worth to note that ReLu neurons can sometimes be pushed into states in which they become inactive for essentially all inputs. In this state, no gradients flow backward through the neuron, and so the neuron becomes stuck in a perpetually inactive state and "dies".

$$f(x) = x^+ = \max(0, x),$$

Figure 8. ReLu activation function.

5 Methods

In order for the agent to interface with the QWOP game environment, it had to be able to simulate keyboard input as well as read the raw pixels on the screen. This was achieved by creating a virtual environment in Python for the Deep Q Learning agent to get the current state and step through actions. Another variable that was needed was the current distance that the runner has traveled. However, due to the obfuscated nature of the native JavaScript game code, we had to rely on other methods to extract the current distance. We utilized the OpenCV library to find image contours of the numbers and corresponding wrapping rectangles. The raw pixels at those locations are then screenshotted, cropped and fed into a support vector machine trained to predict its corresponding number. The Python Imaging Library (PIL) was used to take screenshots of the game and to feed it as raw input into the agent. PyAutoGUI was used to simulate keyboard input.

Since there are four possible inputs into the QWOP game interface, and because buttons can be pressed concurrently, an alternative key schema was defined instead of modeling the actions as four distinct outputs. There are now 16 distinct outputs, each representing a combination of four keys. This schema is defined below in Figure 9. Each row represents one of the 16 possible 4-key combinations and the 1s and 0s respectively represent if that corresponding key is pressed or released.

	Q	W	O	P
A	0	0	0	0
B	0	0	0	1
C	0	0	1	0
D	0	0	1	1
E	0	1	0	0
F	0	1	0	1
G	0	1	1	0
H	0	1	1	1
I	1	0	0	0
J	1	0	0	1
K	1	0	1	0
L	1	0	1	1
M	1	1	0	0
N	1	1	0	1
O	1	1	1	0
P	1	1	1	1

Figure 9. Key input schema definition.

Initially, an environment representing the QWOP game is instantiated. Then an agent is created. For each episode, the agent either steps through predicted actions and receives a reward until it falls over and the game resets, or the agent executes all 500 timesteps. Every tenth episode, the current weight and biases in the neural networks are cached in a backup file. We limit the input to be a small rectangle covering the runner's lower torso and upper thighs in an effort to reduce the time to train the convolutional neural network. The reward is defined by how long the agent stays alive. Thus, the longer the ragdoll runner is alive, the greater the reward will be. The Q-function is incentivized to choose actions that correspond with stability.

6 Results

Empirically, one reliable way to stay alive is to either hold no keys down or press the keys that will result in the runner with its legs spread apart as far as possible. Initial trials with 1000 episodes of 500 timesteps each yielded promising results. The hyperparameters were set as follows in Figure 10. The agent will start off by guessing 100% of its actions and every subsequent episode will decrease the guessing rate by 0.5%. For fear of overshooting, the learning rate was defined to be 0.0001. However, one tradeoff was that it took a significant amount of time for the neural network to converge on the optimal Q-function.

Episodes	1000
Gamma	0.95
Epsilon	1.0
Epsilon decay	0.995
Learning rate	.0001

Figure 10. Hyperparameters for initial trials.

As more episodes were executed, the agent learned to press the same key over and over again. The key combination that found the most success was “J”, which corresponds to holding the "Q" and “P” key down. This configuration allowed the runner to get in a position similar to someone doing the lunges. This position proved to be the most stable, as repeated presses of "Q" and "P" after entering the lunge position is unable make the agent fall over. Due to the low learning rate and low epsilon decay rate, each training session took upwards of eight hours. However, given the hyperparameters, the Deep Q Learning agent learned to start pressing the same keys around episode 300. Then around episode 500, the key combination pressed converged to “J”, providing the most stability to the runner.

With a working Deep Q Learning agent, we attempted to shorten the training time by increasing the learning rate and epsilon decay: the agent guesses less initially and finds global minima faster. However, it is important to note that a very small learning rate causes the network to converge extremely slowly, and if it is too high, we risk overshooting and never finding the global minima. By changing these hyperparameters, the agent was able to learn to press the keys “Q” and “P” repeatedly by episode 200. However, since the agent is staying alive longer, this did not significantly decrease our experimentation time.

A trend was observed between the action variability and episode number. The variability is calculated dividing the most common action count by the total action count. This equation is shown below in Figure 11. A variability value close to zero means that many different combinations are pressed throughout the episode. A variability value close to one means that the same combination was pressed throughout the episode. We can observe that the agent learns that pressing the same buttons tend to result in a higher reward. As the number of training episodes increases, the variability also increases. The plot is shown below in Figure 12.

$$\textit{Variability} = \frac{\textit{Most common action count}}{\textit{Total action count}}$$

Figure 11. Variability equation.

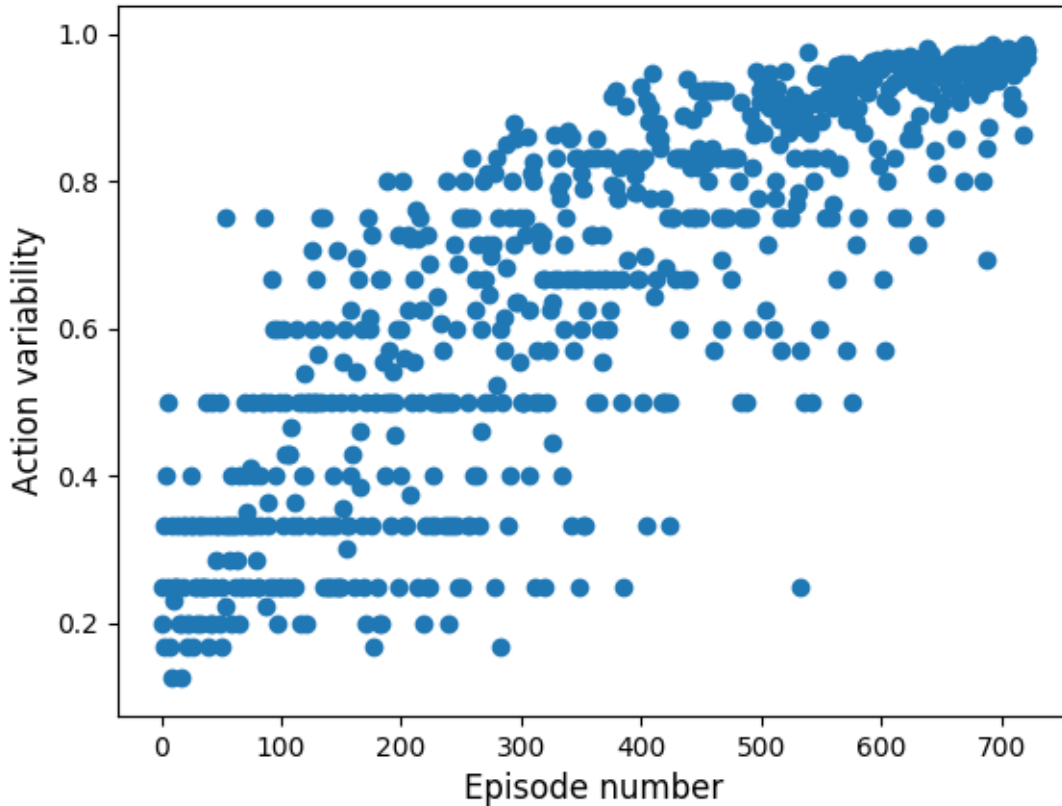


Figure 12. Action variability versus Episode.

A similar trend to the one shown previously in Figure 11 can be observed between the number of actions executed and episode number. This is shown below in Figure 13. We observe that after 500 episodes, the agent was able to stay alive consistently through the 500 timesteps in each episode. Both plots show a slight exponential growth trend, which is expected. As the network learns the correct sequence of actions to take, they are predicted more often and thus result in a higher reward, creating a positive feedback loop.

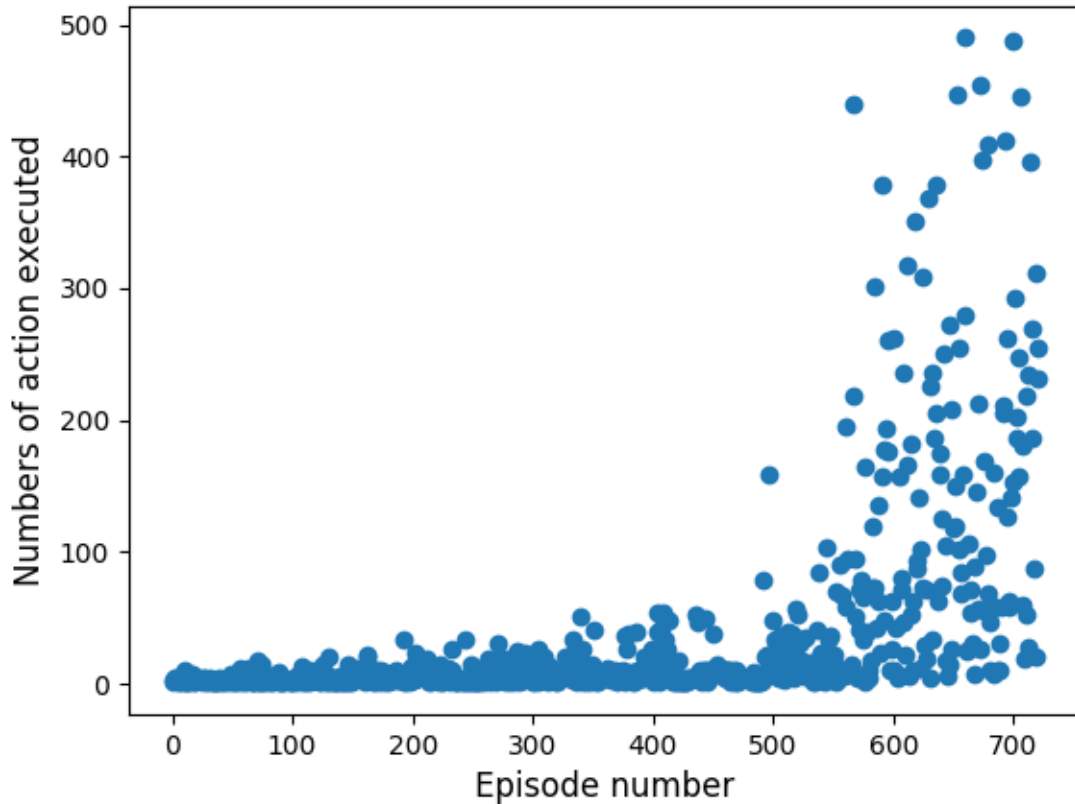


Figure 13. Actions executed versus Episode.

7 Conclusion

7.1 Summary

In this paper, we discussed applying Deep Q Learning to the nonconventional control task of keeping the QWOP runner alive as long as possible. This is in contrast to the traditional way that success is measured in QWOP. Typically, success is defined as distance traveled; however, we redefined the problem and were able to successfully apply our model. We have shown that with only raw pixel inputs, a convolutional neural network can converge to the optimal value of the Q-function. After roughly half of the expected 1000 training episodes, the agent learned to stay alive by holding down the keys “Q” and “P”.

7.2 Future work

Work can be done to modify the current model to play the flash game QWOP as originally intended. Currently, the Deep Q Learning model is incentivized to stay alive for as long as possible. It would be interesting to modify the rewards to incentivize the agent to travel longer distances. Further work can also be done to decrease the latency of OpenCV image processing to find the contours of the distance numbers faster. Faster score detection would mean that there is less delay between consecutive key presses. This model can also be theoretically applied to more complicated environments in OpenAI Gym. Specifically, bipedal and quadrupedal walking environments.

8 References

- [1] Melo F. “Convergence of Q-Learning: a simple proof.”
users.isr.ist.utl.pt/~mtjspaam/readingGroup/ProofQlearning.pdf.
- [2] Wawrzynski P., A. Pacut. “Model-Free off-Policy reinforcement learning in continuous environment.” Proc. IEEE International Joint Conference on Neural Networks, 2004.
- [3] Ng A. “Shaping and Policy Search in Reinforcement Learning.” Ph.D. Dissertation.
University of California, Berkeley. AAI31305322.
- [4] Peters J., Vijayakumar S., Schaal S. “Reinforcement Learning for Humanoid Robotics.”
IEEE-RAS International Conference on Humanoid Robots. 2003.
- [5] Strehl A., Li L., Wiewiora E., Langford J., Littman M. “PAC model-free reinforcement learning.” Proc. 23rd Int'l Conf on Machine learning (ICML), pp 881-888, 2006.
- [6] Mnih V., Kavukcuoglu K., Silver D., Graves A., Antonoglou I., Wierstra D., Riedmiller M.
“Playing Atari with Deep Reinforcement Learning.” arXiv preprint arXiv:1312.5602.
2013.
- [7] Brockman G., Cheung V., Pettersson L., Schneider J., Schulman J., Tang J., Zaremba W.
OpenAI Gym, <https://github.com/openai/gym>, arXiv preprint arXiv:1606.01540. 2016.
- [8] Chollet F. Keras. <https://github.com/fchollet/keras>, GitHub, 2015.
- [9] Brodman G., Voldstad R. “QWOP Learning.” 2012.
- [10] Foddy B. “Foddy.net – Games by Bennett Foddy.” Foddynet, www.foddy.net/.
- [11] Altman E. “Constrained Markov decision processes.” CRC Press. 1999.
- [12] Krizhevsky A., Sutskever I., Hinton G. “ImageNet Classification with Deep Convolutional
Neural Networks.” Proc. 25th Int'l Conf. on Neural Information Processing Systems
(NIPS), pp 1097-1105, 2012.

- [13] Arulkumaran K., Deisenroth M., Brundage M., Bharath A.,” Deep Reinforcement Learning: a brief survey.” IEEE Signal Processing Magazine, Vol 34, N. 6, pp 26-38.
- [14] Goodfellow I., Bengio Y., Courville A. “Deep Learning”, MIT Press, 2016.
- [15] Volodymyr Mnih et al, “Human-level control through deep reinforcement learning”, Nature, 26 February 2015, Vol 51, pp 529-533.