

## Background: Reward functions

The goal of the arm is to approach and consistently touch the cylinder with any part of the arm, so the first reward function gives rewards based on whether the goal was achieved over 100 frames. If the arm touches the ground or the wrong part of the arm, a negative reward is given instead. There are interim rewards based on whether or not the arm is moving towards the objective. One is a reward directly proportional to the movement while the other is based on a slower moving average of the velocity. I also added a small negative interim reward based on goal distance to encourage it to find the object quicker.

The constants REWARD\_WIN and REWARD\_LOSS were set to 1000 and 100 respectively. All of the interim and loss rewards are a ration of REWARD\_LOSS since all rewards in the system except the reward for winning is negative. While the reward for timing out or touching the object with the arm is the full negative loss reward, touching the ground gives is  $\text{REWARD\_LOSS} * \text{goalDistance} / 3$ . This rewards the robot for closing the gap before failing and funnels the actions towards the goal.

The interim reward is:

$$\exp(-2 * \text{avgGoalDelta}) * \text{REWARD\_LOSS} / 50 + \exp(-2 * \text{distDelta}) * \text{REWARD\_LOSS} / 50;$$

The alpha for avgGoalDelta was 0.4 which was slow enough to capture the general direction of the gripper. I also used the instantaneous distance change in the reward inspired by PD controllers. Lastly, I used exponential functions to convert the rewards to a negative based system that exponentially punishes moving away from the goal.

## Background: Hyperparameters

The first parameter tuned was the image input size. They were set at 64x64 to be as low as possible as to still retain the information necessary for the job while optimizing for training speed. Learning rate was initially set to 0.1 since I wanted it to learn faster to hit the 90% accuracy with less episodes. Batch size was set to 64 to speed up learning as well.

I experimented with LSTM and size but it didn't seem to help too much. It gives the benefit of looking at previous frames to decide the next action, but it ends up being more parameters to learn, slowing down the training process. It could be more useful in velocity based joint control.

I used a gamma of 0.95 since I wanted to slightly encourage quicker episodes. But setting it too low would give less weight to the precision at the end of the episode, needed to get the gripper collision. I set the episode to start at .95 randomness and decay to .01. The goal here was to have high randomness in the beginning to encourage exploration and to dial it down once the more efficient paths have been hammered out. I lowered the ending randomness to .01 to reduce the likelihood of failing at the end when high precision is key.

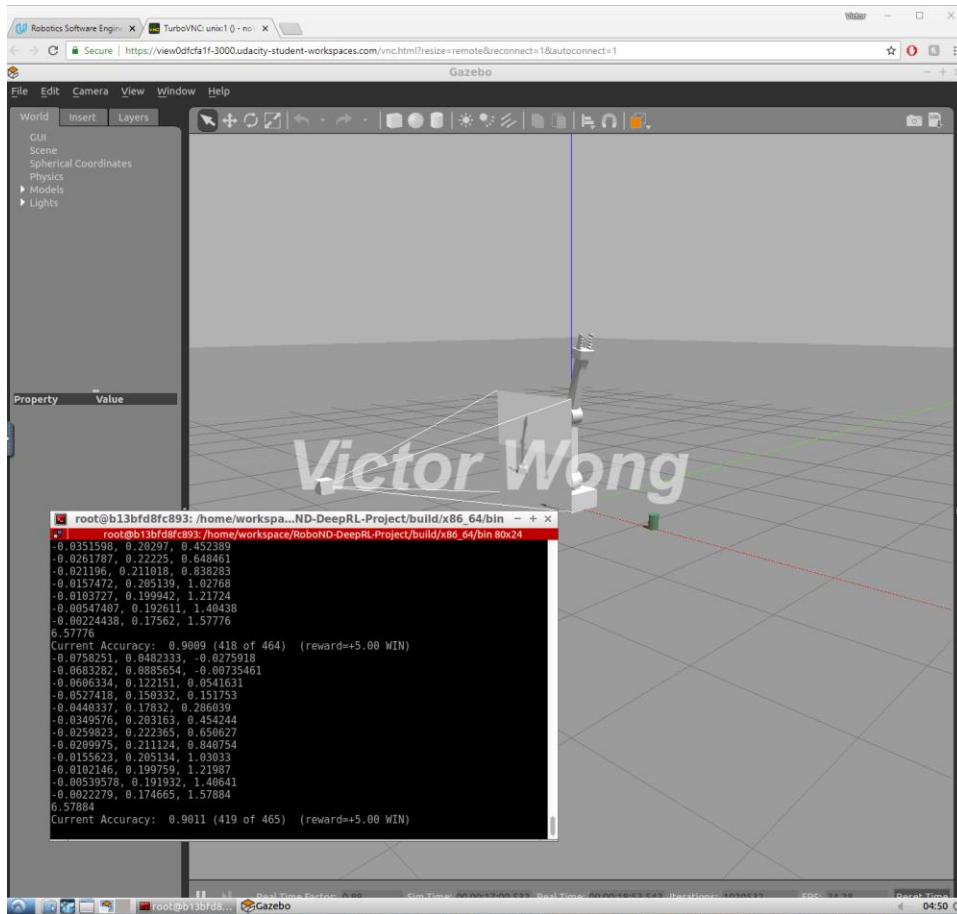
In order to achieve the second task, I had to slightly decrease the learning rate and increase the delay for randomness decay. But most of it boiled down to simplifying my interim reward functions and normalizing their rewards so that no reward would overshadow the other ones. I changed all of the interim reward functions to be negative because I noticed it would get stuck trying to maximize marginal gains from a good move. Especially since the moving average reward takes a couple iterations to decay.

In terms of rewards

## **Results**

I was able to hit just a little over a 90% success rate on the first task and an 82% success rate on the second task. The first one goal was pretty straightforward to hit, since it learned to just learned to apply torque to the first joint and swing the arm against the object. It didn't have to learn how to manipulate the 2<sup>nd</sup> joint so it was easier to get a higher precision.

The second goal took a lot more trial and error to achieve. It boiled down to getting good training sessions in the beginning and locking those actions in. I could have lowered the learning rate and randomness decay even more but it would accrue a lot more failures in the beginning which would take proportionally more wins to offset.





## Future Work

Exploring lower learning rates and longer training times would make the solution more robust. Also adding control to the gripper's fingers or unlocking and restricting the base joint would allow the robot to learn to swing into the object which would be a more robust way of colliding with only the object and not the ground.