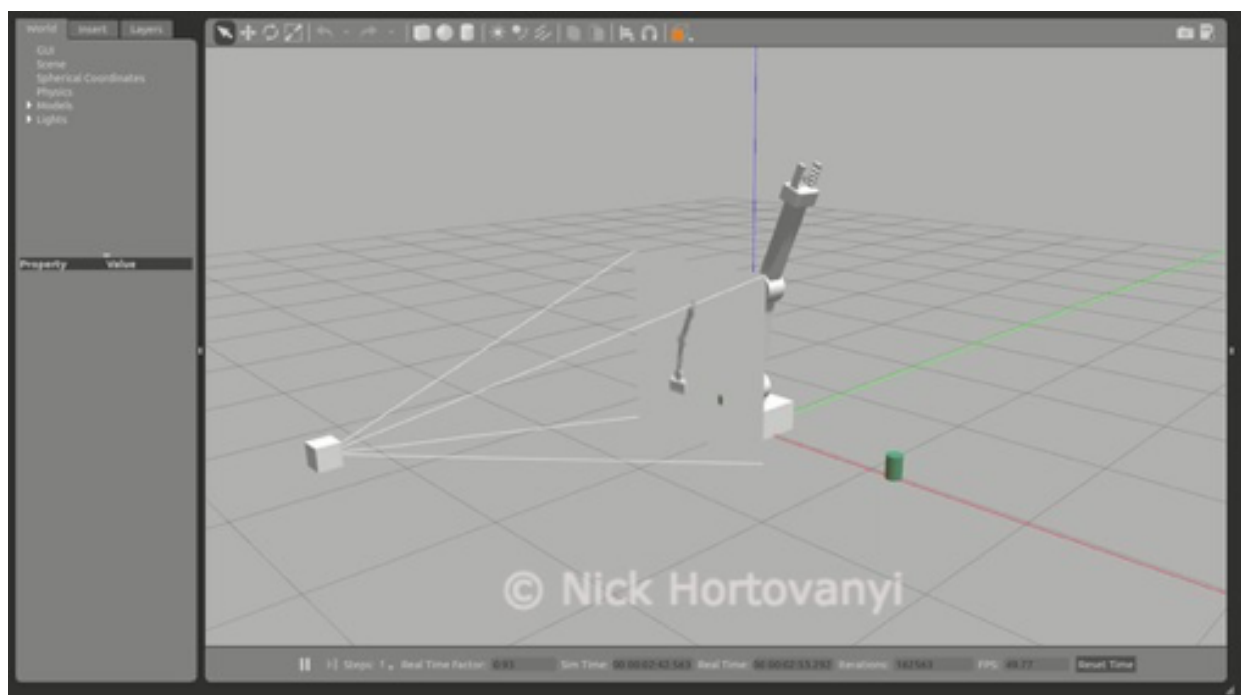


Deep RL Arm Manipulation

This Deep Reinforcement Learning Arm Manipulation project has two objectives to achieve using a template project. Whereby to achieve each objective we create a [DQN](#) agent and define reward functions to teach a robotic arm.

The template project is based on the Nvidia open source project "jetson-reinforcement" developed by [Dustin Franklin](#).



Robot Arm in Gazebo simulator

The two primary project objectives are:

- Have any part of the robot arm touch the object of interest, with at least a 90% accuracy for a minimum of 100 runs.
- Have only the gripper base of the robot arm touch the object, with at least a 80% accuracy for a minimum of 100 runs.

Written by Nick Hortovanyi

Reward functions

The reward functions are defined in [ArmPlugin.cpp](#). The arm joints were updated using position control (as that was the programs default setting). For each joint there are two actions (either increase or decrease joint position).

`REWARD_WIN` was set to `0.125` (`0.1` 2nd objective) with `REWARD_LOSS` to `-0.125` (`-0.1` 2nd objective).

If the robot gripper hit the ground a `REWARD_LOSS * 10` was given and the episode ended.

Interim rewards, within the episode, were issued if there was no ground contact or 100 frames had not been exceeded nor the

The main interim reward was based on the distance goal delta between the gripper and the cycling prop. If a positive weighted average was derived then a `REWARD_WIN` was recorded otherwise `REWARD_LOSS * distance to goal` was issued. Thus the `REWARD_LOSS` was higher the further away from the goal the arm was.

For the gripper base (2nd) objective an additional `REWARD_LOSS` was added if the absolute average goal delta was `< 0.001` to penalise no movement.

If the robot arm hit the prop, a `REWARD_WIN * 10` was used for the first objective otherwise a `REWARD_LOSS * 5` for the second objective if the collision was not with the `gripper_middle`.

However for the second objective a `REWARD_WIN * 20` was issued if the collision point was `gripper_middle`.

Any collision ends the episode.

Hyper Parameters

Image dimensions were set to the same size as the input. Training was

performed on a GTX1070 and there was no need to restrict memory usage.

`INPUT_WIDTH 64` `INPUT_HEIGHT 64`

`OPTIMIZER "Adam"` was chosen as it in general performs better than `RMSPProp` whilst maintaining its advantages.

For objective 1 the `LEARNING_RATE` was `0.1` with `REPLAY_MEMORY` at `1000`. The value was chosen via trial and error.

For objective 2 the `LEARNING_RATE` was decreased to `0.01` due to the higher `REPLAY_MEMORY` set at `20000`. The higher `REPLAY_MEMORY` was used so as to allow for more discrete learning, due to the smaller surface area required to achieve a collision to meet objectives.

For both `BATCH_SIZE` was set to 512 (again sufficient memory on the GTX 1070).

`LSTM` was used `USE_LSTM true` with `LSTM_SIZE 256` which was set via trial and error.

Results

Objective 1 - Have any part of the robot arm touch the object of interest, with at least a 90% accuracy for a minimum of 100 runs.

```
nick@fig: ~/Udacity/RoboND-DeepRL-Project/build/x86_64/bin
Current Accuracy: 0.8750 (049 of 056) (reward=+1.25 WIN)
Current Accuracy: 0.8772 (050 of 057) (reward=+1.25 WIN)
Current Accuracy: 0.8793 (051 of 058) (reward=+1.25 WIN)
Current Accuracy: 0.8814 (052 of 059) (reward=+1.25 WIN)
Current Accuracy: 0.8833 (053 of 060) (reward=+1.25 WIN)
Current Accuracy: 0.8852 (054 of 061) (reward=+1.25 WIN)
Current Accuracy: 0.8871 (055 of 062) (reward=+1.25 WIN)
Current Accuracy: 0.8889 (056 of 063) (reward=+1.25 WIN)
Current Accuracy: 0.8906 (057 of 064) (reward=+1.25 WIN)
Current Accuracy: 0.8923 (058 of 065) (reward=+1.25 WIN)
Current Accuracy: 0.8939 (059 of 066) (reward=+1.25 WIN)
Current Accuracy: 0.8955 (060 of 067) (reward=+1.25 WIN)
Current Accuracy: 0.8971 (061 of 068) (reward=+1.25 WIN)
Current Accuracy: 0.8986 (062 of 069) (reward=+1.25 WIN)
Current Accuracy: 0.9000 (063 of 070) (reward=+1.25 WIN)
Current Accuracy: 0.9014 (064 of 071) (reward=+1.25 WIN)
Current Accuracy: 0.9028 (065 of 072) (reward=+1.25 WIN)
Current Accuracy: 0.9041 (066 of 073) (reward=+1.25 WIN)
Current Accuracy: 0.9054 (067 of 074) (reward=+1.25 WIN)
Current Accuracy: 0.9067 (068 of 075) (reward=+1.25 WIN)
Current Accuracy: 0.9079 (069 of 076) (reward=+1.25 WIN)
Current Accuracy: 0.9091 (070 of 077) (reward=+1.25 WIN)
Current Accuracy: 0.9103 (071 of 078) (reward=+1.25 WIN)
```

© Nick Hortovanyi

Results objective 1

The robotic arm quickly learnt how to hit the prop with a degree accuracy in a repeatable fashion. On occasion if the arm trained initially away from the prop, it would take longer to achieve a higher accuracy.

Once a winning path was learnt this configuration consistently had the robotic arm quickly hitting the prop objective.

As can be seen in the above summary output the objective was achieved well within the criteria specified.

Objective 2 - Have only the gripper base of the robot arm touch the object, with at least a 80% accuracy for a minimum of 100 runs.

```
nick@fig: ~/Udacity/RoboND-DeepRL-Project/build/x86_64/bin
Current Accuracy: 0.7816 (068 of 087) (reward=+2.00 WIN)
Current Accuracy: 0.7841 (069 of 088) (reward=+2.00 WIN)
Current Accuracy: 0.7865 (070 of 089) (reward=+2.00 WIN)
Current Accuracy: 0.7778 (070 of 090) (reward=-1.00 LOSS)
Current Accuracy: 0.7802 (071 of 091) (reward=+2.00 WIN)
Current Accuracy: 0.7826 (072 of 092) (reward=+2.00 WIN)
Current Accuracy: 0.7849 (073 of 093) (reward=+2.00 WIN)
Current Accuracy: 0.7872 (074 of 094) (reward=+2.00 WIN)
Current Accuracy: 0.7895 (075 of 095) (reward=+2.00 WIN)
Current Accuracy: 0.7917 (076 of 096) (reward=+2.00 WIN)
Current Accuracy: 0.7938 (077 of 097) (reward=+2.00 WIN)
Current Accuracy: 0.7959 (078 of 098) (reward=+2.00 WIN)
Current Accuracy: 0.7980 (079 of 099) (reward=+2.00 WIN)
Current Accuracy: 0.8000 (080 of 100) (reward=+2.00 WIN)
Current Accuracy: 0.8020 (081 of 101) (reward=+2.00 WIN)
Current Accuracy: 0.8039 (082 of 102) (reward=+2.00 WIN)
Current Accuracy: 0.8058 (083 of 103) (reward=+2.00 WIN)
Current Accuracy: 0.8077 (084 of 104) (reward=+2.00 WIN)
Current Accuracy: 0.8095 (085 of 105) (reward=+2.00 WIN)
Current Accuracy: 0.8113 (086 of 106) (reward=+2.00 WIN)
Current Accuracy: 0.8131 (087 of 107) (reward=+2.00 WIN)
Current Accuracy: 0.8148 (088 of 108) (reward=+2.00 WIN)
```

© Nick Hortovanyi

Results objective 2

With the finer control required, and alteration to the interim reward system, this configuration would often hesitate before making a move. Whilst it learnt quickly how to get very close to having the `gripper_middle` hit the prop, it would also often just miss either hitting the ground or the arm itself hitting the prop. There seemed to be a repeatable pattern, of just extending past and swinging down in an arch, that once learnt gave consistent winning results.

Occasionally the middle joint would collide with the ground and this would lead to the objective not being met.

This configuration was not always reproducible, however with the above screen shot it was able to meet the objectives.

Future work

There were clear arcs that once found achieved a win quickly. Such that it would be worthwhile investigating an interim reward system based on not just the distance from the goal but also distance from an ideal arc trajectory as the arm approached.

Further using centre points to calculate distance from goals becomes less accurate the closer to the goal the arm is. Such that other points like the end of the gripper_middle and top of prop cylinder, would be worthwhile experimenting with.