

Project 7: Deep RL Manipulator Arm

Ahmed Khalid

May 9, 2019

1 Reward Function

The reward function is the same for the 2 objectives:

- +100 if the final objective is achieved
- -100 if the the episode exceeds the maximum number of frames, the arm hits the ground or hits itself
- Interim reward is calculated as follows: $\text{reward} = 10 * (\text{smoothed moving average}) - 0.5$

The smoothed moving average is calculated for the distance between the gripper and the can. While the 0.5 is a time penalty that is introduced to encourage the agent to achieve the objective faster and not settle to collecting smaller rewards without actually reaching the can.

Position control was used to solve both of the objectives. The agent was much faster to converge using the position control. Although the movement of the arm was too quick and not realistic in the simulation. Speed control was much slower to converge although the movements of the arm were smoother.

2 Hyperparameters

The hyperparameters were chosen as follows:

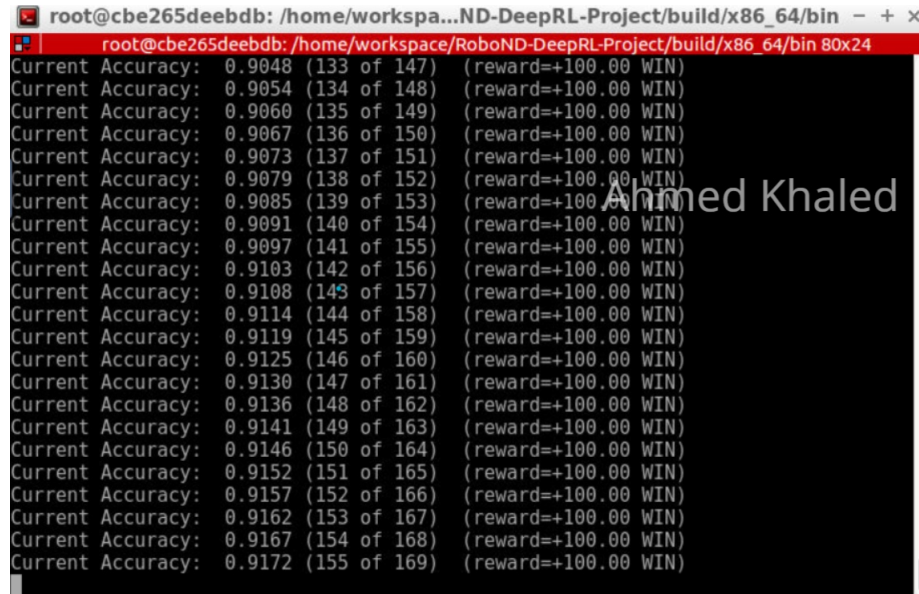
- INPUT_WIDTH 64 & INPUT_HEIGHT 64: The dimensions of the input image is set according to 2 criteria, performance and accuracy. If it is too small the details will vanish from the image and decrease the accuracy. If it is too high is will take longer to execute and require more parameters to train. The value of 64 achieved a good balance between performance and accuracy.
- OPTIMIZER "Adam": It was chosen because in general performs better than RMSProp whilst maintaining its advantages.

- **LEARNING_RATE 0.01:** It was chosen based on experimentation. In general, a low value of learning rate causes slower but less noisy learning. While a higher value may make the learning faster but it may not converge to an optimal minimum of the loss function.
- **BATCH_SIZE 64:** It was chosen based on trial and error. The higher the batch size the more steady the learning gets but a value too high caused the agent to take longer to converge.
- **REPLAY_MEMORY 10000:** This value is also chosen based on trial and error. The replay memory is used to keep batch gradient updates stable and to enable the agent to learn from past experiences.
- **USE_LSTM true & USE_LSTM 64:** The LSTM was used to keep track of the past experiences of the agent and not take actions based only on the current input or frame. The value of USE_LSTM was based on trial and error.

3 Results

3.1 Objective #1: Have any part of the robot arm touch the can

The agent converged after about 20 episodes. Because the task is relatively easy. It achieved an accuracy of 90% after about 130 episodes.



```

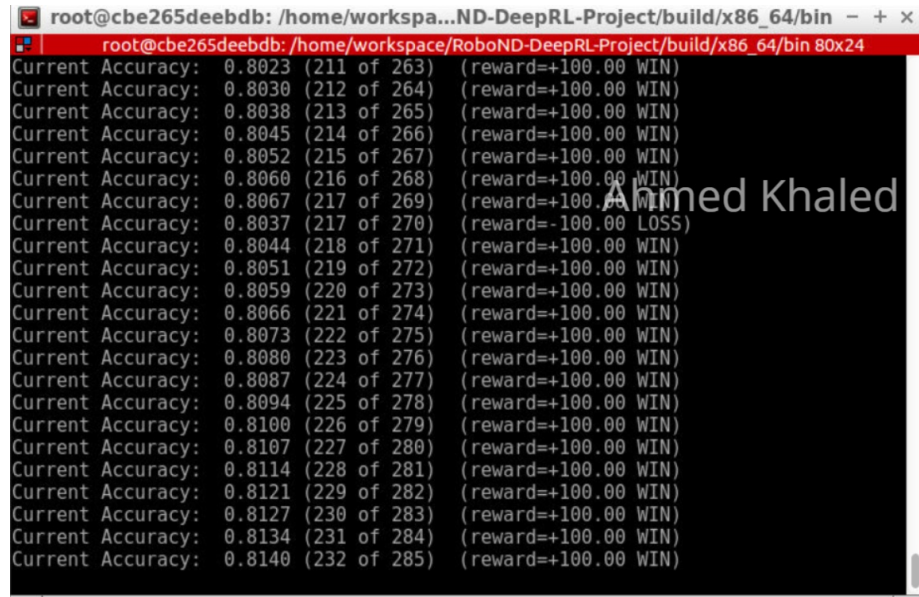
root@cbe265deebdb: /home/workspa...ND-DeepRL-Project/build/x86_64/bin - + x
root@cbe265deebdb: /home/workspace/RoboND-DeepRL-Project/build/x86_64/bin 80x24
Current Accuracy: 0.9048 (133 of 147) (reward=+100.00 WIN)
Current Accuracy: 0.9054 (134 of 148) (reward=+100.00 WIN)
Current Accuracy: 0.9060 (135 of 149) (reward=+100.00 WIN)
Current Accuracy: 0.9067 (136 of 150) (reward=+100.00 WIN)
Current Accuracy: 0.9073 (137 of 151) (reward=+100.00 WIN)
Current Accuracy: 0.9079 (138 of 152) (reward=+100.00 WIN)
Current Accuracy: 0.9085 (139 of 153) (reward=+100.00 WIN)
Current Accuracy: 0.9091 (140 of 154) (reward=+100.00 WIN)
Current Accuracy: 0.9097 (141 of 155) (reward=+100.00 WIN)
Current Accuracy: 0.9103 (142 of 156) (reward=+100.00 WIN)
Current Accuracy: 0.9108 (143 of 157) (reward=+100.00 WIN)
Current Accuracy: 0.9114 (144 of 158) (reward=+100.00 WIN)
Current Accuracy: 0.9119 (145 of 159) (reward=+100.00 WIN)
Current Accuracy: 0.9125 (146 of 160) (reward=+100.00 WIN)
Current Accuracy: 0.9130 (147 of 161) (reward=+100.00 WIN)
Current Accuracy: 0.9136 (148 of 162) (reward=+100.00 WIN)
Current Accuracy: 0.9141 (149 of 163) (reward=+100.00 WIN)
Current Accuracy: 0.9146 (150 of 164) (reward=+100.00 WIN)
Current Accuracy: 0.9152 (151 of 165) (reward=+100.00 WIN)
Current Accuracy: 0.9157 (152 of 166) (reward=+100.00 WIN)
Current Accuracy: 0.9162 (153 of 167) (reward=+100.00 WIN)
Current Accuracy: 0.9167 (154 of 168) (reward=+100.00 WIN)
Current Accuracy: 0.9172 (155 of 169) (reward=+100.00 WIN)

```

Figure 1: Terminal screenshot of the first objective

3.2 Objective #2: Have only the gripper base of the robot arm touch can

The agent converged after about 50 episodes. Because the task is relatively harder than the first one. It achieved an accuracy of 80% after about 200 episodes. It should be noted that increasing the time penalty too much caused the arm to just slam into the can without actually caring about achieving the original objective.



```
root@cbe265deebdb: /home/workspa...ND-DeepRL-Project/build/x86_64/bin - + x
root@cbe265deebdb: /home/workspace/RoboND-DeepRL-Project/build/x86_64/bin 80x24
Current Accuracy: 0.8023 (211 of 263) (reward==+100.00 WIN)
Current Accuracy: 0.8030 (212 of 264) (reward==+100.00 WIN)
Current Accuracy: 0.8038 (213 of 265) (reward==+100.00 WIN)
Current Accuracy: 0.8045 (214 of 266) (reward==+100.00 WIN)
Current Accuracy: 0.8052 (215 of 267) (reward==+100.00 WIN)
Current Accuracy: 0.8060 (216 of 268) (reward==+100.00 WIN)
Current Accuracy: 0.8067 (217 of 269) (reward==+100.00 WIN)
Current Accuracy: 0.8037 (217 of 270) (reward=-100.00 LOSS)
Current Accuracy: 0.8044 (218 of 271) (reward==+100.00 WIN)
Current Accuracy: 0.8051 (219 of 272) (reward==+100.00 WIN)
Current Accuracy: 0.8059 (220 of 273) (reward==+100.00 WIN)
Current Accuracy: 0.8066 (221 of 274) (reward==+100.00 WIN)
Current Accuracy: 0.8073 (222 of 275) (reward==+100.00 WIN)
Current Accuracy: 0.8080 (223 of 276) (reward==+100.00 WIN)
Current Accuracy: 0.8087 (224 of 277) (reward==+100.00 WIN)
Current Accuracy: 0.8094 (225 of 278) (reward==+100.00 WIN)
Current Accuracy: 0.8100 (226 of 279) (reward==+100.00 WIN)
Current Accuracy: 0.8107 (227 of 280) (reward==+100.00 WIN)
Current Accuracy: 0.8114 (228 of 281) (reward==+100.00 WIN)
Current Accuracy: 0.8121 (229 of 282) (reward==+100.00 WIN)
Current Accuracy: 0.8127 (230 of 283) (reward==+100.00 WIN)
Current Accuracy: 0.8134 (231 of 284) (reward==+100.00 WIN)
Current Accuracy: 0.8140 (232 of 285) (reward==+100.00 WIN)
```

Figure 2: Terminal screenshot of the first objective

4 Future Work

As the current implementation of the reward function has no regard to the actual trajectory that the arm takes. It may be beneficial to experiment on modifying it to prioritize an optimal trajectory. Also implementing more advanced reward functions from research papers could be really beneficial. It can also benefit from using better training techniques, like learning rate annealing with restarts. The annealing part is used to make the network make more steady progress as it reaches a minimum. The restarts would also make the network get out of a local minimum if it gets stuck into one.