# Project: Deep RL Arm Manipulation

Oleksii Zhurbytskyi

**Abstract**—The goal of the project is to create a DQN agent and define reward functions to teach a robotic arm to carry out two primary objectives:
1. Have any part of the robot arm touch the object of interest, with at least a 90% accuracy.
2. Have only the gripper base of the robot arm touch the object, with at least a 80% accuracy.

**Index Terms**—Robot, IEEEtran, Udacity, LATEX, Deep RL.

---✦---

## 1 REWARD FUNCTIONS

### 1.1 Objective 1: Any part of the robot arm touch the object of interest

1) If any part of robot touch the object the reward is REWARD_WIN * 10
2) If any part of robot touch the ground the penalty is REWARD_LOSS * 10
3) Additional reward/penalty if the robot do not touch the ground depends of the smoothed moving average of the delta of the distance to the goal

### 1.2 Objective 2: Only the gripper base of the robot arm touch the object

1) If the gripper base of robot touch the object the reward is REWARD_WIN * 20
2) If any other part of robot touch the object the penalty is REWARD_LOSS * 5
3) If any part of robot touch the ground the penalty is REWARD_LOSS * 10
4) Additional reward/penalty if the robot do not touch the ground depends of the smoothed moving average of the delta of the distance to the goal

The joint movements were controlled by position.

## 2 HYPERPARAMETERS

The main factors in choosing hyperparameters were the size of the input image, the amount of RAM and learning results. Two optimizers were tested, and it turned out that the Adam shows better results than the RMSProp for this task. LEARNING_RATE and REPLAY_MEMORY hyperparameters have been selected by trial and error starting from 0.1 for LEARNING_RATE and 1000 for REPLAY_MEMORY. If the goal could not be achieved, the LEARNING_RATE was decreased and the REPLAY_MEMORY was increased.

- INPUT_WIDTH is 64 and INPUT_HEIGHT is 64 - were chosen according to the size of the image
- OPTIMIZER is Adam
- LEARNING_RATE is 0.1 for Objective 1 and 0.01 for Objective 2
- REPLAY_MEMORY is 1000 for Objective 1 and 10000 for Objective 2

- BATCH_SIZE is 256 as we have a large amount for RAM
- LSTM_SIZE is 256

Thus, the only differences in hyperparameters for two objectives are LEARNING_RATE and REPLAY_MEMORY.

## 3 RESULTS

For both objectives specified accuracy has been achieved in the provided environment. It turned out that the sooner the robot reaches the goal for the first time, the faster the training goes. The provided screenshots show that the goals were achieved within 100 runs.
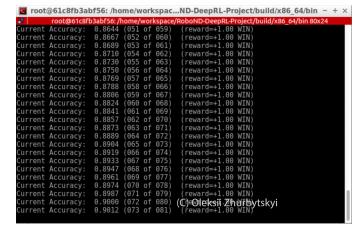


Fig. 1. Result for Objective 1 - any part of the robot arm touch the object of interest, with at least a 90% accuracy.

## 4 FUTURE WORK

As a result of the experiments, it became obvious that training begins to go well only after the robot first touches the target. Therefore, further improvements should be concentrated precisely on this task. Most likely this can be achieved by optimizing the reward functions.

Fig. 2. Result for Objective 2 - only the gripper base of the robot arm touch the object, with at least a 80% accuracy.