# Deep RL: Manipulator

Shilpaj Bhalerao

July 29, 2019

**Abstract:** *Deep Reinforcement Learning has achieved a great deal of achievements. AlphaGo used Deep RL to win game against world's best Go player. This paper shows use of Deep RL by a robotic arm manipulator to touch the object.*

**Keywords:** ROS, Gazebo, Deep Learning, Reinforcement Learning, Deep RL, etc.

## 1. Introduction

Deep Reinforcement learning is a field of robotics in which raw image is taken from a camera and fed to Deep Neural Network. The output of the DNN is passed to an agent which gets trained on the data and then applies control value to the actuators. In this paper, a camera is placed in front of the robotic arm. The camera input is passed to the DRL agent.
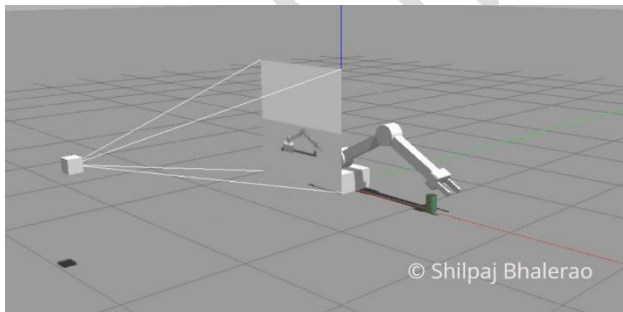


Fig 1: Robotic ARM manipulator with camera

The DRL agent has to achieve following tasks:

- 90% accuracy when arm collides with prop
- 80% accuracy when gripper collides with prop

## 2. Reward Function

The reward function is defined in ArmPlugin.cpp in the repository. The objective of the project is to touch the prop. The episode will end if the robotic arm touches the ground or the object.

To achieve this task, the arm is punished when it touches the ground or it fails to complete the task within the selected frames.

The reward function is designed to reward the arm 50 points when it touches the object and penalizes it 50 points if it fails to complete the task or touches the ground.

An interim reward was given based on the current position of the robotic arm. The method used is a smooth moving average based on delta value which is a distance of arm from the object. The reward function also contains alpha term which decides how much to rely on the previous delta value. This reward function eliminates the fluctuations and helps to maximize the reward.

## 3. Hyperparameters

```
#define INPUT_WIDTH    64
#define INPUT_HEIGHT   64
#define OPTIMIZER "Adam"
#define LEARNING_RATE 0.1f
#define REPLAY_MEMORY 10000
#define BATCH_SIZE 256
#define USE_LSTM true
#define LSTM_SIZE 256
  ©Shilpaj Bhalerao
#define REWARD_WIN  50.0f
#define REWARD_LOSS -50.0f
#define INTERIM_REWARD 4.0f
#define INTERIM_OFFSET 0.3f
#define ALPHA 0.4f
```

Fig 2: Hyperparameters

Above figure shows the hyperparameters used for the project. The input size is reduced from 256 to 64 in order to reduce the processing time. Adam optimizer is used for the DRL network. Usually, 0.01 is good initial value of learning rate but in this case, it wasn't helping so decided to use 0.1. Batch size 256 is used along with the LSTM of 256 units. The value of alpha which decides how to rely of previous value is 40%.

## 4. Results

Objective-1 of achieving task with more than 90% accuracy when arm touches the prop is tried with two different reward values of 20 and 50. The objective is achieved earlier in case of 20 points compared to the 50 points.

Fig 3: Objective-1 with more than 90% accuracy

Fig 4: Objective-1 with different reward function

Fig 5: Objective-2 with more than 80% accuracy

As the agents keeps training for a greater number of epochs, it keeps on improving its performance. Initially, it couldn't even figure out in which direction to move. As the time passes, it becomes more and more consistent.

To achieve the goal of more than 90% accuracy, agent is rewarded if there is a collision between arm and the prop. Similarly, in case of 80% accuracy, agent is only rewarded when gripper collides with prop.

The agent has learned to perform its objective but more fine tuning on system with higher processing capability will help to improve the results.

## 5. Future Work

To improve the project results, one can try using a deeper neural network. The network can be optimized using Batch Normalization, dropout, etc. One can also try making the network wider to achieve more accuracy.

I want to try placing 2 cameras perpendicular to each other in the gazebo environment which will allow one more degree of freedom to the robot. Then the robot can be trained to pick up objects in 3D world.