

Deep Reinforcement Learning for Autonomous Systems

Piero Macaluso - s252894
Candidate
Politecnico di Torino

Prof. Pietro Michiardi
Supervisor
EURECOM

Prof. Elena Baralis
Supervisor
Politecnico di Torino

This document represents the summary of the master thesis project. The source code of this work is publicly available at <https://github.com/pieromacaluso/Deep-RL-Autonomous-Systems>

1. Introduction

Because of its potential to thoroughly change mobility and transport, autonomous systems and self-driving vehicles are attracting much attention from both the research community and industry. Recent work has demonstrated that it is possible to rely on a comprehensive understanding of the immediate environment while following simple high-level directions, to obtain a more scalable approach that can make autonomous driving a ubiquitous technology. However, to date, the majority of the methods concentrates on deterministic control optimisation algorithms to select the right action, while the usage of deep learning and machine learning is entirely dedicated to object detection and recognition.

Recently, we have witnessed a remarkable increase in interest in Reinforcement Learning (RL). It is a machine learning field focused on solving Markov Decision Processes (MDP), where an agent learns to make decisions by mapping situations and actions according to the information it gathers from the surrounding environment and from the reward it receives, trying to maximise it. As researchers discovered, it can be surprisingly useful to solve tasks in simulated environments like games and computer games, and it showed encouraging performance in tasks with robotic manipulators. Furthermore, the great fervour produced by the widespread exploitation of deep learning opened the doors to function approximation with convolutional neural networks, developing what is nowadays known as deep reinforcement learning.

1.1. Objective

In this Thesis, we argue that the generality of reinforcement learning makes it a useful framework where to apply autonomous driving to inject artificial intelligence not only in the detection component but also in the decision-making

one. The focus of the majority of reinforcement learning projects is on a simulated environment. However, a more challenging approach of reinforcement learning consists of the application of this type of algorithms in the real world.

After an initial phase where we studied the state-of-the-art literature about reinforcement learning and analysed the set of possible alternatives about technologies to use, we started our project starting from the ideas contained in [2], where the authors were able to train a self-driving vehicle by using Deep Deterministic Policy Gradient (DDPG) [3] by tuning hyper-parameters in simulation. We decided to not use simulators in our approach, therefore we researched an algorithm suitable for real-world experiments and capable of work fine without an expensive real-world hyper-parameter tuning. We found in Soft Actor-Critic (SAC) [1] the algorithm we needed.

Therefore, our thesis consisted of two main contribution:

1. Designing of the Control System to let all components and technologies involved interact;
2. Setting up the experimental framework with SAC algorithm and carrying out an experiment to analyse strengths and weaknesses of this approach.

2. Design of the control system

We based our project on Cozmo, a little toy robot produced by Anki, whose developers offered a granular and fully-featured Python SDK with many interfaces to allow a direct control of the robot. We found it to be suitable for our experimental needs: it mounts a grayscale camera with 60 FOV and has two tracks to steer and drive. Our aim was to apply deep reinforcement learning algorithms, so we decided to use PyTorch as deep learning framework and the standardised approach provided by OpenAI Gym to build up the reinforcement learning environment for the experiment.

Our idea was to build up a system where the human has the total control on the experiment flow to directly teach to the robot how to drive: he is able to start the episode and to stop it when the robot reaches a pernicious situation in

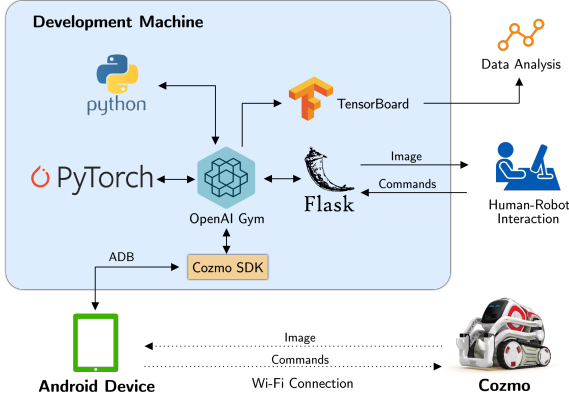


Figure 1. Outline of the control system that shows the most crucial technologies and component involved.

order to reposition it in the closest safe situation and restart the loop. In this scenario, the human has total responsibility for how the robot is trained: he is the one who decides when an action is dangerous or not by disengaging reinforcement learning algorithms decision.

To obtain this configuration, we managed to design a simple and intuitive user interface that prompts the user when he started an experiment. This interface provides a live stream from Cozmo on-board camera and a set of key with the related function. It works through Javascript to communicate to a Flask server that communicates directly to Cozmo SDK and the OpenAI Gym environment to provide information for the user (e.g. images, learning information) and the robot (e.g. commands). We used TensorBoard to gather data for posterior analysis and representations. The system we obtained is outlined in figure 1.

In order to start the experiment, we formalised the problem as a MDP obtaining the following setup:

State Space The current state observation consists of two subsequent 64×64 grayscale images which represent the situation before and after the action respectively. We opted for resizing the images provided by Cozmo SDK for memory consumption reason due to the large experience replay memory needed for off-policy algorithms.

Action Space The representation of how the robot can interact with the surrounding environment is given by a vector of two real value. The first value represents the desired speed and its range is between 0 and 1, while the second one represents the steering of the vehicle and its range is between -1 and 1. Both values are the result of a simplification useful in the learning phase. In practise, the system designed converts these values to properly interact with the robot. The maximum speed reachable is 150mm/s.

Reward Function The reward function defines what is the ultimate objective of the task formalised. After an analysis of the literature available, we decided to formalise this

function as the total length of track crossed by the robot after each action taken. This decision revealed its simplicity, but also its effectiveness. Furthermore, this choice leads to match reward with length traveled, useful factor to allow the developer to easily quantify the agent improvements and have a better feedback.

CozmoDriver-v0 is the Cozmo environment we designed. It represents the result of a continuous development and testing to solve all the problems and particularities that we have found along the way. We implemented all the requirements and needs dictated by the type of experiment, putting in place all the compromises to manage the available RAM memory and make experiments manageable by reducing the learning time between one episode and the next one. One of the most crucial features we implemented is the rescue system, useful in a context where experiments last tens of hours and unexpected fails can occur. It consists of two main part: the first is a volatile backup system for every single episode that is discarded after the start of the next one and allows the user to make the agent forget a faulty episode, while the second one is a checkpoint system to save the state of the experiment and allow the user to restore it in the following days.

3. Experiments

The path that led to the final implementation of the algorithm allowed us to detect some specific requirements and to overcome them with appropriate countermeasures. We decided to firstly implement a simplified environment to test functionalities and reinforcement learning algorithms we aimed to use in Cozmo environment. We used the inverted pendulum swing-up problem, a classic problem in the control literature and available in OpenAI Gym. The original implementation of this environment consisted of observations with values related to the current angle and speed of the pendulum. For this reason we decided to build a wrapper for the original environment in order to receive observations as raw pixels, since the goal was to apply the same considerations and the same convolutional neural networks that we would use in the Cozmo environment.

The results of the experiments carried out using hyper-parameters taken from the available literature, showed that SAC algorithm has a better performance than the DDPG one, both in terms of stability and number of episodes to achieve a working policy. Therefore, we decided to implement SAC algorithm to carry out experiments in the Cozmo environment.

We opted for a neural network with three convolutional layers with 16 features of 3×3 dimension, a stride of 2, zero padding with two fully-connected layers with a hidden size of 256 features in the last part. We applied batch normalisation after each convolutional level and Rectified Linear Unit function (ReLU) as non-linearity.

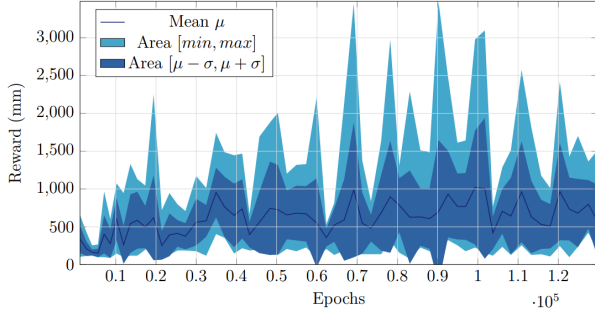


Figure 2. SAC CozmoDriver-v0 Test Average Reward Plot. The graph reports mean, standard deviation range and min-max range of the average reward obtained from 10 test episodes every 50 episodes.

We managed to complete a whole set of 3000 episodes exploiting the SAC algorithm to solve the autonomous driving task with Cozmo. Taking into account waiting times between episodes and charging times, we managed to complete the experiment in almost one working week, after almost 1.3×10^5 epochs of learning.

The agent reached the most crucial results in the testing phase presented in figure 2 where we plotted minimum and maximum values obtained in every set of ten episode together with the mean and the standard deviation. Following this approach we noticed a performance increase with a maximum mean reached of almost 1 metre. Furthermore, the maximum value reached among all tests episodes was equal to almost 3.5 metres which equals more than one complete tour of the track. It is noticeable that the results are not stable as we expected from the experiments with the inverted pendulum environment: the reward values do not improve uniformly with increasing epochs. However, carrying out the experiments episode by episode, we noticed a marked improvement in the performance obtained in the tests. The robot learned to approach turns and to stay on the lane of a straight road.

4. Conclusions

The plot reports a visible improvement in the maximum length of track traveled before the disengagement of the user. Despite these improvements, the agent was not able to learn how to drive in a secure and stable way, as we can notice from the unstable growth of the mean reward. These facts made us reflect on the critic points of our experiment setup that may have had a role in the instability of the results obtained.

We localised two major problems which, in our opinion, have had a particular influence on the results obtained. The first factor was the amount of RAM available in the development machine. This limitation forced us to decrease the size of the replay memory and a consequent early deletion

of less recent episodes. Analysing the plots, we noticed that this fact translated in the increase of the temperature parameter: this symptom underlines the need for the algorithm to explore more the solution space. The second major problem was the limitation of the camera sensor on the robot, particularly its viewing angle. The features offered by the Anki Cozmo camera proved to be inadequate to observe the track we designed. We noticed this fact after many episodes when the robot started to improve its performance: it began to adopt a wave behaviour on the straights, interpreting the vision of a single road line as a curve.

4.1. Future Work

Our proposals about future improvements to the project grow from the weakness in our approach. It could be interesting to execute these algorithms on a device with a bigger RAM, but also to design this approach with a Variational Auto-Encoder (VAE) to reduce the dimensionality of the information retrieved during experiments.

It may be useful to enhance sensors installed in the self-driving robot. A possible alternative to Anki Cozmo could be Anki Vector, the successor of Cozmo which mounts a 720p camera with 120 Ultra Wide FOV or to build up a personal *Donkey Car* with custom specifications. Anki Vector could be interesting to perform reinforcement learning algorithms with the usage of the renewed front camera together with the infrared laser scanner on-board to investigate approaches to data fusion with more sensors.

Another intriguing research path consists of an investigation about the application of model-based reinforcement learning algorithms to autonomous driving. A more in-depth review of the literature to better understand the feasibility of this approach, focusing on its strengths and weaknesses compared to model-free ones can be the right starting point.

References

- [1] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- [2] A. Kendall, J. Hawke, D. Janz, P. Mazur, D. Reda, J.-M. Allen, V.-D. Lam, A. Bewley, and A. Shah. Learning to drive in a day. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8248–8254. IEEE, 2019.
- [3] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.