

Deep Reinforcement Learning with Continuous Control in CARLA

Moritz Zanger¹Florian Rottach¹Izel Kilinic¹

Abstract—Style:

- Written last
- 100-250 words
- Past tense

Content:

- Condensed version of entire article

I. INTRODUCTION

The task of teaching vehicles how to drive autonomously in urban scenarios is a challenging and complex one to solve. Not only is there the problem of finding the adequate response to a given situation but also the challenge of taking into account the surrounding factors that have an influence on the state that a vehicle is in and its possible actions. To date, most approaches focus on the manual design of behavioral policies, such as defining a driving policy through the use of annotated maps. While these solutions might work in situations which are documented by the provided mapping infrastructure, they are often difficult to generalize or scale, as they do not necessarily enable the comprehension of any given local scene. In order to make autonomous driving truly feasible in a real-world scenario it would be better to develop systems which are able to find their way without having to rely on an explicit set of rules. One possible solution to this task is provided by reinforcement learning methods. Here, the agent, i.e. the vehicle, actively searches for the optimal driving policy whilst trying to maximize a numerical reward signal. As opposed to imitation learning techniques, which have been popular in finding driving policies (1), reinforcement learning algorithms enable a car to exceed human abilities, if applied correctly. In recent years, deep reinforcement learning methods have proven to be successful in solving complex tasks such as playing GO (2) or Atari (3) and there have been efforts in tackling various problems in the field of autonomous driving, including continuous control tasks (4).

However, two major drawbacks of reinforcement learning methods are their heavy dependency on adequate input state representations (5) and, as with other machine learning techniques, their need of a sufficient amount of accurate sample data to train on. In order to be able to train safely on an adequate amount of data, one approach is the use of data from other domains. For this purpose, the urban driving simulator CARLA has been developed, which is used as a simulation environment for this project.

In this paper, several state-of-the-art reinforcement learning algorithms are implemented and compared, with regard

to their performance considering different driving tasks. Additionally, reward functions for the respective problems are tested and several input representations are designed and evaluated.

?? What is it exactly that we contributed that is new to already existing research??

II. RELATED WORK

III. BACKGROUND

We regard the typical reinforcement learning setting where an agent interacts with an environment \mathcal{E} . At each one of a number of discrete timesteps t the agent decides on taking an action a_t from a given set of actions \mathcal{A} . This is done based on the state s_t that the agent is currently in and following a policy π , which is a mapping of the possible states to the action space $\pi : \mathcal{S} \rightarrow \mathcal{P}(\mathcal{A})$. In this case, π is stochastic, as it returns the probability distribution over the possible actions, and the action space is continuous.

As a result, the agent receives a reward r_t and the subsequent state s_{t+1} of his environment. The setup is assumed to follow the properties of a Markov decision process, where besides state space \mathcal{S} , action space \mathcal{A} and reward function $r(s_t, a_t)$, we also include the transition function to the future states $p(s_{t+1}|s_t, a_t)$.

IV. CONCEPT/METHODS AND MODELS

In the course of this project, three different algorithms are used. Following the work of Lillicrap et al. (4) and Mnih et al. (6), the (DDPG)- and (A3C) algorithms are implemented and compared according to their performance in the Gym environment CarRacing-v0. Later on, the PPO algorithm is applied to solve a continuous control task in the CARLA simulator. TODO: Hier unser Konzept analog zur Präsentation vorstellen (Slide mit den 4 Kästen)

A. DDPG

One very notable advance in reinforcement learning has been made by the development of the so-called "Deep Q Network" (Mnih et al., 2015). The DQN is able to solve tasks with high-dimensional observation spaces. However, it is only efficiently capable of working with discrete and low-dimensional action spaces. In order to adapt a (DQN) for the successful use with continuous control problems, as given in CarRacing-v0, a discretization of the action space has to be carried out, which can lead to two main difficulties: an explosion in the number of possible actions and the loss of important information (4).

To evade these obstacles (Lillicrap et al.) propose a new approach, called the Deep Deterministic Policy Gradient

¹The authors are with FZI Research Center for Information Technology, Haid-und-Neu-Str. 10-14, 76131 Karlsruhe, Germany {kuhnt, zoellner}@fzi.de

(DDPG), which is a model-free, off-policy actor-critic algorithm. They adopt the advantages of (DQN) and combine them with the actor-critic framework, resulting in the stabilization of Q-learning by using a replay buffer and soft updates on the target networks of both actor and critic, through

$$\tau \ll 1 : \theta' \leftarrow \tau \theta + (1 - \tau) \theta' \tag{1}$$

and finding a deterministic policy.
 Noise with Ornstein-Uhlenbeck

- B. A3C
- C. PPO
- D. Reward function

The design of the reward function turned out to be a much more difficult procedure as expected initially. The already provided rewards in gym’s CarRacing environment were rather simple and still led to a good result in the end. For CARLA however, it was necessary to invest more time into the engineering and designing of a suitable inducement system, because the driving behavior depended on more factors and the performance metric was not just driving as fast as possible. In the following, the process of arriving at our final reward function will be described. In the first step, we investigated possible sensor inputs that might have an impact on the driving behavior and discussed on their impact. The outcome is summarized in the following table:

TODO:	Schön	formatieren
Per frame penalty	Forces agent to move	
Lane invasion counter	As few lane changes as possible	
Steering angle change	Avoids oszillations	
Delta heading	Drive as straight as possible relative to road	
Position change	Maximize travelled distance	
Collision binary per frame	Avoid crashes	
Velocity	Drive fast under the other constraints	
Distance to middle lane	Drive as centered as possible	

We started with incrementally adding these attributes to our reward calculation and quickly realized, that the main challenge is to adjust the weights and harmonize the contrary effects of the terms. An example would be, that giving the velocity a relatively high weight, such as 0.8, while giving the distance to center line a weight of 0.2 results in an agent that speeds over the map and pays few attention on lane invasions. On the opposite side, the agent will drive only very slowly or even not at all, if the rewards for oscillations are too high compared to the velocity. Considering, that we have not only two, but rather several possible components, this results in a complex combinatorical problem that can either be solved by trail and error or applying permutation optimization techniques. To achieve initial results, we started to discover a proper reward function ”by hand”. We pruned the above list by applying the following considerations that resulted from tesing different approaches.

Firstly, some components show a redundant behavior and hence one of them can be removed. An example would be the velocity and the position change - both contain the same information. Further, the attributes lane invasion and steering angle didnn’t affect the driving behavior in a positive way. The two most important parameters turned out be the velocity and the delta heading, which expressed the relative angle to the current street angle. Our final reward function and the aggregation weights can be found in the following table:

TODO:	Werte	aus	Präsentation	übernehmen
	hier	drei	spalten	

The result is very promising and can be summarized as follows: The agent is capable of driving smoothly within the right lane and can perform turn manuevers on most intersections. It attempts to drive around other vehicles, but is not capable of breaking. We trained different models on this reward function and all of them had a good performance compared to other functions.

- E. Input representation
 - TODO: Analog zu Präsentation aufbauen
- F. Autoencoder/Latent space
 - TODO: Encoder-decoder architecture bzw. Generator nennen
- G. Training
 - TODO: Training beschreiben analog zur präsentation (vgl folien)

V. EVALUATION

- A. CarRacing-v0
- B. CarRacing-v0
- C. Results
- D. Comparison algorithms/reward functions

VI. CONCLUSIONS

- Similar to the abstract but more detail
- Conclusion of the key points of each section
- Summary of main findings
- Important conclusions that can be drawn
- Discuss benefits and shortcomings of our approach
- Suggest future areas of research

REFERENCES