



Controlling Self-Driving Race Cars with Deep Neural Networks

UNIVERSITY OF OSNABRÜCK

DEPARTMENT OF NEUROINFORMATICS

BACHELOR'S THESIS

Author:
Christoph Stenkamp

Supervisors:
Prof. Dr. Gordon Pipa
Leon Sütfeld

Osnabrück,
August 11, 2017

Abstract

This Thesis will be written in the next two months, and I'm pretty scared about that.

Preface

This document was written as the author's bachelor thesis at the department of neuroinformatics at the University of Osnabrück during summer 2017 and is an original and independent work by the author Christoph Stenkamp.

Christoph Stenkamp
Osnabrück, August 11, 2017

Acknowledgements

Thanks to my parents, Marie, my supervisors, and my friends...

“There are no surprising facts, only models that are surprised by facts; and if a model is surprised by the facts, it is no credit to that model.”

Eliezer Yudkowsky

Contents

Abstract	i
Preface	ii
Acknowledgements	iii
1 Introduction	1
1.1 Motivation	1
1.1.1 Problem Domain	1
1.1.2 Goal of this thesis	1
1.2 Research Questions	1
1.3 Reading Guidelines	1
2 Reinforcement Learning	2
2.1 Reinforcement Learning Problems	2
Markov Decision Processes	2
Value of a state	3
Value of an action	4
2.2 Temporal difference Learning	5
SARSA	5
Q-learning	5
2.3 Q-Learning with Neural Networks	6
2.3.1 Deep Q-learning	6
Double-Q-Larning	7
Dueling Q-Learning	7
2.3.2 Deterministic Policy Gradient	7
2.3.3 Exploration techniques	7
3 Related work	8
3.1 Reinforcement Learning Frameworks	8
3.2 Self-driving cars	8
4 Program Architecture	9
4.1 Design choices	9
4.1.1 The game as a reinforcement learning problem	9
4.1.2 The vectors	9
4.1.3 Exploration	9
4.1.4 Reward	9
4.1.5 Performance measure	9
4.2 Implementation	9
4.2.1 The game	9
What Leon did already	9
Communication	9

4.2.2	The agent	9
	Pretraining	9
	The different agents	9
	Challenges and Solutions	9
5	Analysis, Results and open Questions	10
6	Discussion	11
7	Conclusion and future directions	12
	Declaration of Authorship	13

List of Figures

List of Tables

List of Algorithms

List of Abbreviations

The abbreviations used throughout the work are compiled in the following list below. Note that the abbreviations denote the singular form of the abbreviated words. Whenever the plural forms is needed, an s is added. Thus, for example, whereas ANN abbreviates *artificial neural network*, the abbreviation of *artificial neural networks* is written ANNs.

ANN	Artificial Neural Network
CNN	Convolutional (artificial) Neural Network
CPU	Central Processing Unit
DDPG	Deep Deterministic Policy Gradient - Network
DQN	Deep-Q-Network
GUI	Graphical User Interface

List of Symbols

For my friends, family, and especially Marie.

Chapter 1

Introduction

1.1 Motivation

1.1.1 Problem Domain

1.1.2 Goal of this thesis

1.2 Research Questions

1.3 Reading Guidelines

Chapter 2

Reinforcement Learning

As the task at hand was not only to provide a reinforcement learning agent, but also to convert a game itself into something the agent can successfully play, I will in this chapter go into detail about Reinforcement Learning in general, to give insights into why I did what I did. Also, I will try to keep this stuff as general as possible, getting into detail when speaking about the used algorithms. [The sense of this chapter is to give an intro of MDPs and RL. It shall also go into enough details on how to specify an MDP such that an RL agent can learn on it, because a big part of the work was exactly that. It's supposed to end with SARSA and Q-learning as the two Ideas on how to perform RL]

2.1 Reinforcement Learning Problems

Machine Learning can mainly be subdivided into three main categories: Supervised Learning, Unsupervised Learning, and Semi-supervised learning. The first deals with direct classification or regression using labelled data (i.e. it uses pairs of data-points with their corresponding category or value). In unsupervised learning, no such label exists, and the data must be clustered into meaningful parts without any knowledge, by for example grouping objects by similarity of their properties. What will be mainly considered in this thesis will be a certain kind of semi-supervised learning: *Reinforcement learning*. In Reinforcement Learning (**RL**), instead of labels for the data, there is a *weak teacher*, which provides feedback to the actions the agent took.

Markov Decision Processes

The metaphor behind RL is that of a decision maker (*agent*) and an *environment*. The agent makes observations in the environment (its input), takes actions (output) and receives rewards. In contrast to the classical ML approaches, in RL the agent is also responsible for exploration, as he has to acquire his knowledge actively. Thus, a reinforcement learning problem is given if the only way to collect information about the *underlying model* (the environment) is by interacting with it. As the environment does not explicitly provide actions the agent has to perform, its goal is to figure out the actions maximizing its cumulative reward until a training episode ends.

In the classical RL approach, the environment is divided into discrete time steps. If that is the case, the environment corresponds to a *Markov Decision Process* (**MDP**). Formally, a MDP is a 5-tuple $\langle S, A, P, R, \gamma \rangle$, consisting of the following:

S – set of states $s \in S$

A – set of actions $a \in A$

$P_a(s, s')$ – transition probability function from state s to state s' under action a

$R_a(s, s')$ – reward function for action a in state s if the environment moves to s'

γ – discount factor for future rewards $0 \leq \gamma \leq 1$

Though in general both the state transition function and the reward function may be indeterministic (and thus not in complete control of the decision maker), I will refer to the result of a state transition at discrete point in time t as $s_{t+1} := P(s_t, a_t)$ and to the result of the reward function as $r_t := R(s_t, a_t, s_{t+1})$. If no point in time is explicitly specified, it is assumed that all variables use the same t .

While an *offline learner* gets as input the problem definition in the form of a complete MDP, where the only task left is to classify actions yielding high rewards from actions giving suboptimal results, the task for an *online reinforcement learning* agent is a lot harder, as it has to learn the MDP itself via trial and error. In the process of reinforcement learning, the agent will encounter states s of the environment, performing actions a . The future state s_{t+1} of the environment may be indeterministic, but depends on the history of previous states s_0, \dots, s_t as well as the action of the agent a_t . It is assumed that the *Markov property* holds, which means that a state s_{t+1} depends only on the current state s_t and current action a_t .

Throughout interacting with the environment, the agent receives rewards r , depending on his action a as well as the state of the environment s . In many RL problems, the full state of the environment is not known to the agent, and it only perceives an observation depending on the environment: $o_t := o(s_t)$ ¹. This is referred to as *partial observability*, and the corresponding decision process is a *partially observable MDP*. Additionally, the agent knows when a final state of the environment is reached, terminating the current training episode. An episode consists for the agent thus of a sequence of observations, actions and rewards until at time t_t some terminal state s_{t_t} is reached:

$$\text{Episode} := ((s_0, a_0, r_0), (s_1, a_1, r_1), (s_2, a_2, r_2), \dots, (s_{t_t}, a_{t_t}, r_{t_t}))$$

Value of a state

In the process of reinforcement learning, the agent tries to perform as well as possible in the previously unknown environment. For that, it uses an action-policy π , mapping states to actions: $\pi(s) = a$.² In general, this policy may also be stochastic. As the agent does not have supervised data for what actions are the ground truth, it must learn some kind of measure for the value of being in a certain state or performing a certain action. The commonly used measure for the value of a state can be calculated by the immediate reward this state gives, summed with the discounted future reward the agent will achieve by continuing to follow his policy from this

¹From now on, when I mean the state of the environment, I will explicitly refer to it as s_e , while reserving s for the agent's observation of the environment $o(s_e)$

²It is obvious, that the result of both the reward function and the state transition function depend on π . To be explicit about that, I will refer to a reward dependent on π as r^π and a state transition dependent on π as s^π . If state or reward depends on an explicit action instead, I refer to it as r^a and s^a .

state on:

$$V^\pi(s_t) := \sum_{t'=t}^{t_t} (\gamma^{t'-t} * r_{t'}^\pi) \quad (2.1)$$

Using the discounted future reward is useful because in an indeterministic environment it gets less likely that the agent actually reaches this state, and to make the agents perform good actions as quickly as possible.

The actual, underlying Value of a state s is defined as the value of the state when using the best possible policy, which corresponds to the maximally achievable reward starting in state s :

$$V^*(s_t) := \max_\pi V^\pi(s_t^\pi) \quad (2.2)$$

While *passive reinforcement learning* simply tries to learn the Value-function without the need of action selection, an *active reinforcement learner* tries to estimate a good policy, using which those high-value states are actually reached. If the value of every state is known, then the optimal policy can be defined as the one achieving maximal value for every upcoming state: $\pi^* := \operatorname{argmax}_\pi V^\pi(s) \forall s \in S$. Knowing what an optimal policy does, and using 2.1 and 2.2, it is possible to re-write the definition of the value of a state recursively as

$$V^*(s_t) = \max_\pi \left(\sum_{t'=t}^{t_t} (\gamma^{t'-t} * r_{t'}^\pi) \right) \quad (2.3)$$

$$= \max_\pi (r_t^\pi + \gamma * V^\pi(s_{t+1}^\pi)) \quad (2.4)$$

This is known as the *Bellman Equation*, which allowed for the birth of dynamic programming. It rewrites the value of the decision problem at time t in terms of the immediate reward at t plus the value of the remaining decision problem at $t + 1$, resulting from the initial choices.³

Value of an action

While the definition of a state-value is useful, it alone does not help an agent to perform optimally, as neither the successor function $P_a(s, s')$, nor the reward function $R_a(s)$ are known to the agent. While so-called *model-based* reinforcement learning tries to learn both of those explicitly to reconstruct the entire MDP, *model-free* agents use a different measure of quality: the *Q-value*. It represents the value of performing action a_t in a state s_t , afterwards following the policy π .

$$Q^\pi(s_t, a_t) := r_t^{a_t} + \gamma * V^\pi(s_{t+1}^{a_t}) \quad (2.5)$$

With the optimal, maximally archivable action-value Q^* being respectively

$$Q^*(s_t, a_t) = r_t^{a_t} + \gamma * V^*(s_{t+1}^{a_t}) \quad (2.6)$$

$$= \max_\pi (r_t^{a_t} + \gamma * V^\pi(s_{t+1}^{a_t})) \quad (2.7)$$

³This is because of the definition of Bellman's *Principle of Optimality*, which states that "An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision[Bellman1957]"

As the Value of a state is defined as the maximally achievable reward from that state, the relation between Q and V can be expressed as

$$V(s_t) = \max_{a_t} Q(s_t, a_t) \quad (2.8)$$

When an agent knows the Q -value for each action of a state, it can easily infer the optimal action in state s_t as $a_t^* := \operatorname{argmax}_{a_t} (Q(s_t, a_t))$ and thus the optimal policy π^* , guaranteeing maximum future reward at every state. The goal of a model-free RL agent is thus to get a maximally precise estimate of Q^* , yielding maximal reward for every state. For that, it does not need to explicitly learn the reward- and transition function, but instead can model only the Q -function. Its policy is then to simply always take the action yielding the highest value for every state (a *greedy* policy). In RL settings with a highly limited amount of discrete states and actions, the respective Q -function estimate can be specified as a lookup table, whereas for areas of interest, the function is calculated using a kind of nonlinear function approximator.

Throughout exploration of the environment, the agent collects more information of it, continually updating its estimate Q^π . For that, it uses samples from its episodes of interacting with the environment.

2.2 Temporal difference Learning

Throughout the process of reinforcement learning, the agent continually improves its estimates \hat{Q} of Q^* . An optimal solution would be to calculate the Loss of the current Q -estimate as the squared difference $(\hat{Q} - Q^*)^2$, to perform gradient descent in order to minimize that difference. However, as Q^* is unknown to the agent, that is not possible. Instead, a Q -learning agent performs *iterative approximation*, using the information about the environment, to continually update its estimates of Q^* . Using the recursive definition of a state-value given in the Bellman equation 2.4 allows for a technique called *temporal difference learning* [sutton1988]: At time $t+1$, the agent can compare its estimate of the Q -function of the last step, $\hat{Q}^\pi(s_t, a_t)$, with a new estimate using the new information it gained from the environment: r_{t+1} and s_{t+1} . Because of the newly gained information from the underlying model, the new estimate will be closer to the actual function Q^* than the original value:

$$Q^*(s_t, a_t) = r_t^{a_t} + \gamma * V^*(s_{t+1}^{a_t}) \quad (2.9)$$

$$\approx r_t^{a_t} + \gamma * V^\pi(s_{t+1}^{a_t}) = r_t^{a_t} + \gamma * \hat{Q}^\pi(s_{t+1}^{a_t}, a_{t+1}) \quad (2.10)$$

SARSA

The new knowledge about the environment can be incorporated in two different ways. For the first method, the agent samples a full tuple of $\langle s_t, a_t, r_t, s_{t+1}, a_{t+1} \rangle$ from the environment, to then calculate the temporal difference error in non-terminal states as $TD := (r_t + \gamma * \hat{Q}_i^\pi(s_{t+1}, a_{t+1})) - \hat{Q}_i^\pi(s_t, a_t)$. This algorithm of calculating the temporal difference error is known as SARSA, and it is an example of *on-policy* learning. In on-policy learning, the agent uses his own policy in every estimate of the Q -value.

Q-learning

In contrast to SARSA stands the *off-policy* algorithm *Q-learning*. This algorithm does not need to sample the action a_{t+1} , as it calculates the Q -update at iteration i using

the best possible action in state s_{t+1} ⁴. As the previous definition of Q-values was only correct in non-terminal states, a case differentiation must be introduced for terminals and non-terminal states. In the following, y_t will stand for the updated estimate of the Q-value at t , sampling the necessary states, rewards and actions from the environment, resulting in the formula found in [mnih_human-level_2015]:

$$y_t = \begin{cases} r_t & \text{if } t = t_t \\ r_t + \gamma * \max_{a_{t+1}} (\hat{Q}^\pi(s_{t+1}, a_{t+1})) & \text{otherwise} \end{cases} \quad (2.11)$$

The temporal difference error is accordingly defined as

$$TD_t := y_t - \hat{Q}^\pi(s_t, a_t) \quad (2.12)$$

Using this error straight away allows for the update-rule of an agent in a very limited setting: Consider an agent, specifying his approximation of the Q-function (his *model*) with a lookup-table, initialized to all zeros. It is proven that for finite-state Markovian problems with nonnegative rewards the update-rule for the Q-table $Q(s_t, a_t) \leftarrow r_t^{a_t} + \gamma * Q^\pi(s_{t+1}, a_{t+1})$ converges to the optimal Q^* -function, making the greedy policy π^* optimal⁵. It is noteworthy, that each new temporal difference will affect not only the last prediction, but all previous predictions (learning propagates backwards).

As however in practice the problems using a table as the Q-functions are only very limited scenarios, an update rule like this is irrelevant. Instead, a better idea is to use this definition of the temporal difference error for a loss function, which is to be minimized throughout the process of RL. A commonly used loss-function is the *L2-Loss*, which allows gradient descent, updating the parameters of the Q-function into the direction of the newly acquired knowledge. The L2-Loss at iteration i with model-parameters θ_i is thus defined as the following:

$$L_i(\theta_i) := \left(y_i(\theta_i) - Q_i^\pi(s_t, a_t; \theta_i) \right)^2 \quad (2.13)$$

2.3 Q-Learning with Neural Networks

To understand this section, basic knowledge on how *Artificial Neural Networks* (ANNs) work and what they do is presupposed. As mentioned before, it is not only possible to use a Q-table to estimate the Q^* -function, but any kind of function approximator. Thanks to the universality theorem⁶, it is known that ANNs are an example of such.

2.3.1 Deep Q-learning

When using ANNs as function approximators for the model of the environment, it will result in a Loss function depending on the Neural Network parameters, specified by θ . The update rule in Deep Networks depends on the gradient with respect to the weights of this formula, $\Delta_{\theta_i} L(\theta_i)$. While there are attempts to use Artificial Neural Networks for Q-learning as early as 1993[SOURCE], those had only minor success when being applied to more than toy-problems. But everything changed when Deepmind attacked.

⁴A slight deviation from this *double-Q-learning*, an algorithm I will go into detail about lateron.

⁵Of course the agent will need some kind of exploration technique first, more on that later

⁶<http://neuralnetworksanddeeplearning.com/chap4.html>, I need a better source on this!

Double-Q-Larning

Dueling Q-Learning

2.3.2 Deterministic Policy Gradient

2.3.3 Exploration techniques

Chapter 3

Related work

3.1 Reinforcement Learning Frameworks

Gym/Universe Torcs

3.2 Self-driving cars

Nvidias deep-drive RRT* Tensorkart Tesla Lidar

Chapter 4

Program Architecture

The program was written by the author of this work and is licensed under the GNU General Public License (GNU GPLv3). Its source code is attached in the appendix of this work and additionally can be found digitally on the enclosed CD-ROM. The machine learning part was written in PYTHON, using the TENSORFLOW-library [abadi_tensorflow:_2015].

4.1 Design choices

4.1.1 The game as a reinforcement learning problem

4.1.2 The vectors

4.1.3 Exploration

4.1.4 Reward

4.1.5 Performance measure

4.2 Implementation

4.2.1 The game

What Leon did already

Communication

4.2.2 The agent

Pretraining

The different agents

Challenges and Solutions

hier auch Batchnorm, DQN vs DDPG, sehend vs nicht-sehend, ...

Chapter 5

Analysis, Results and open Questions

testing took place on a win10 machine, ... Answer all of the research questions explicitly!!!

Chapter 6

Discussion

Chapter 7

Conclusion and future directions

Declaration of Authorship

I, Christoph Stenkamp, hereby certify that the work presented here is, to the best of my knowledge and belief, original and the result of my own investigations, except as acknowledged, and has not been submitted, either in part or whole, for a degree at this or any other university.

signature

city, date