



UNIVERSITEIT•STELLENBOSCH•UNIVERSITY  
jou kennisvennoot • your knowledge partner

# **Learning to solve Sliding Puzzles using Reinforcement Learning**

Umr Barends  
18199313

Report submitted in partial fulfillment of the requirements of the module  
Project (E) 448 for the degree Baccalaureus in Engineering in the Department of  
Electrical and Electronic Engineering at Stellenbosch University.

Supervisor: JC Schoeman

October 11, 2020

# Acknowledgements

I would like to thank my dog, Muffin. I also would like to thank the inventor of the incubator; without him/her, I would not be here. Finally, I would like to thank Dr Herman Kamper for this amazing report template.



UNIVERSITEIT • STELLENBOSCH • UNIVERSITY  
jou kennisvennoot • your knowledge partner

## Plagiaatverklaring / *Plagiarism Declaration*

1. Plagiaat is die oorneem en gebruik van die idees, materiaal en ander intellektuele eiendom van ander persone asof dit jou eie werk is.

*Plagiarism is the use of ideas, material and other intellectual property of another's work and to present is as my own.*

2. Ek erken dat die pleeg van plagiaat 'n strafbare oortreding is aangesien dit 'n vorm van diefstal is.

*I agree that plagiarism is a punishable offence because it constitutes theft.*

3. Ek verstaan ook dat direkte vertalings plagiaat is.


*I also understand that direct translations are plagiarism.*

4. Dienooreenkomstig is alle aanhalings en bydraes vanuit enige bron (ingesluit die internet) volledig verwys (erken). Ek erken dat die woordelike aanhaal van teks sonder aanhalingstekens (selfs al word die bron volledig erken) plagiaat is.

*Accordingly all quotations and contributions from any source whatsoever (including the internet) have been cited fully. I understand that the reproduction of text without quotation marks (even when the source is cited) is plagiarism*

5. Ek verklaar dat die werk in hierdie skryfstuk vervat, behalwe waar anders aangedui, my eie oorspronklike werk is en dat ek dit nie vantevore in die geheel of gedeeltelik ingehandig het vir bepunting in hierdie module/werkstuk of 'n ander module/werkstuk nie.

*I declare that the work contained in this assignment, except where otherwise stated, is my original work and that I have not previously (in its entirety or in part) submitted it for grading in this module/assignment or another module/assignment.*

Studentenommer / 18199313	Handtekening: 
Initials and surname / U.Barends	Datum / October 11, 2020

# Abstract

## English

The English abstract.

# Contents

<b>Declaration</b>	<b>ii</b>
<b>Abstract</b>	<b>iii</b>
<b>List of Figures</b>	<b>v</b>
<b>List of Tables</b>	<b>vi</b>
<b>Nomenclature</b>	<b>vii</b>
<b>1. Introduction</b>	<b>1</b>
1.1. Section heading . . . . .	1
<b>2. Solvability of a NxN sliding puzzle</b>	<b>2</b>
2.1. General puzzle description . . . . .	2
2.2. Conditions for solvability . . . . .	2
<b>3. Reinforcement Learning Theory</b>	<b>4</b>
<b>4. Summary and Conclusion</b>	<b>6</b>
<b>Bibliography</b>	<b>7</b>
<b>A. Project Planning Schedule</b>	<b>8</b>
<b>B. Outcomes Compliance</b>	<b>9</b>
<b>C. Student and Supervisor agreement</b>	<b>10</b>

# List of Figures

2.1. Example of a sliding puzzle . . . . .	2
2.2. Odd boards with change in blank piece only having even inversion change [1]	3
2.3. Even board solvability [1] . . . . .	3
C.1. Student and Supervisor agreement . . . . .	10

# List of Tables

# Nomenclature

## Variables and functions

$p(x)$                       Probability density function with respect to variable  $x$ .

## Acronyms and abbreviations

AE	Afrikaans English
RL	Reinforcement Learning
MDP	Markov Decision Process



# Chapter 1

## Introduction

Generally in robotics a manipulator (eg. an arm) is used to manipulate an object in the environment. It is normally easier to first simulate the robots behavior in a more simple environment. In this project we try to solve a sliding puzzle using reinforcement learning, where the same algorithm can then later be applied to the robotics problem.

### 1.1. Section heading

# Chapter 2

## Solvability of a NxN sliding puzzle

### 2.1. General puzzle description

Let us assume that we have an NxN puzzle, then we have NxN number of blocks. We can represent the puzzle as an NxN array, then we stack the array into a one dimensional array of  $1 \times (N \times N)$ . For example see the 4x4 puzzle in Figure 2.1 we have a  $1 \times 16$  array as:  $\text{Array} = (12, 7, 8, 13, 4, 9, 2, 11, 3, 6, 15, 14, 5, 1, 10)$ . Before we describe the conditions for a sliding puzzle to be solvable, we first define the term “inversion”. Assuming the the first index of the  $1 \times N^2$  array starts at the left top corner (valued 12) in Figure 2.1, and that it runs from  $[0, (N \times N) - 1]$ . Then an inversion occurs when  $\text{Array}[\text{index}] > \text{Array}[\text{index} + 1]$  where index is an arbitrary integer between 0 and  $N \times N - 1$ . Hence in Figure 2.1 we have a total:  $\text{sum of inversions}(\text{Array}) = 11 + 6 + 6 + 8 + 3 + 5 + 1 + 5 + 1 + 2 + 4 + 3 + 1 + 0 = 56$ .

12	7	8	13
4	9	2	11
3	6	15	14
5	1		10

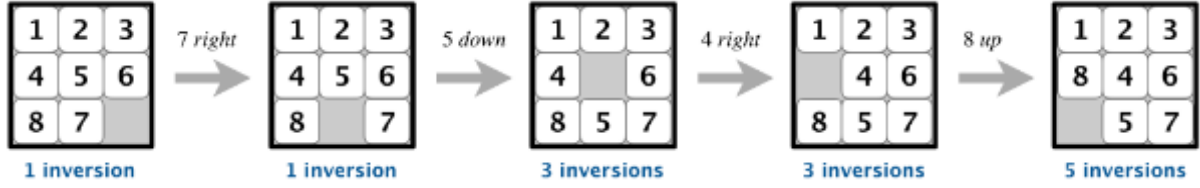
Figure 2.1: Example of a sliding puzzle

### 2.2. Conditions for solvability

Even and odd sized boards are analysed separately (where  $\text{size} = N$ ).

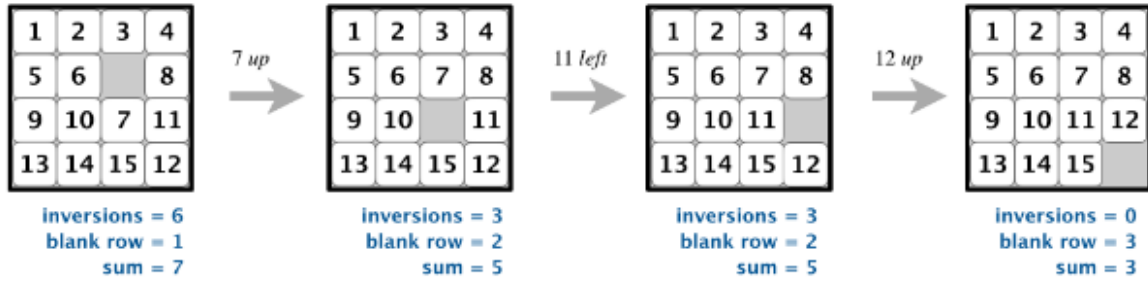
For odd sized boards where  $N$  is odd we have the puzzle only being solvable if and only if the boards has an even number of inversions. The proof for this can be deduced by looking at Figure 2 and noting that for every switch of the blank block we have an even change in the sum of inversions of the board. [1]

For even sized boards where  $N$  is even we have the board solvable if and only if the



**Figure 2.2:** Odd boards with change in blank piece only having even inversion change [1]

number of inversions plus the row of the blank square is odd. This is illustrated in Figure 3.



**Figure 2.3:** Even board solvability [1]

Half of all puzzle configurations are unsolvable. [2] This means that we only have  $N! / 2$  configurations that are solvable for an  $N \times N$  board. This was proven using parity in the paper in [2]. Sliding puzzles can be solved relatively quickly with today's processing of computers for puzzles for example an  $5 \times 5$  puzzle was solved in 205 tile moves in 2016. [3]

The issue more so lies in finding the shortest path to solving a puzzle. This specific problem of solving with the least amount of tile moves of a sliding puzzle has been defined as NP (non-deterministic polynomial-time) hard. NP hardness is are problems that are at least as hard as NP. Where in computational complexity theory NP (non-deterministic polynomial-time) is a has a solution with a proof variable to be in polynomial time by a deterministic Turing Machine. A Turing machine is a mathematical model defining an abstract machine which manipulates symbols according to a set of rules. [4]

In simpler terms a problem is NP if it can be solved within a time that is a polynomial function of the input. For instance if we define the time to solve a problem as 'T' and the input data as 'D'. Then as long as  $T = \text{polynomial function}(D)$  then a problem is NP.

# Chapter 3

## Reinforcement Learning Theory

Reinforcement learning is a method of learning what to do by linking states to actions with a numerical reward incurred for every state-action pair. The final goal of RL is to find a path of states and actions with the maximum amount of reward. [5]

To model reinforcement learning problems we use dynamic systems theory, specifically using incompletely-known MDP's (Markov decision processes). Most of RL problems can be described using MDP's. In Mathematics a MDP is a stochastic, discrete time control process. What that means is that the process is essentially partially controlled and partially random. MDP's depend mainly on a few variables which are, states, actions, state transition probability, reward and a discount factor. These variables can be denoted in a 5-tuple as:

$$(S, A, P_{ss'}^a, R_s^a, \gamma)$$

where

- S is a finite set of states
- A is a finite set of actions
- $P_{ss'}^a$  is a matrix of probabilities with  $P_{ss'}^a = P(S_{t+1} = s' | S_t = s, A_t = a)$
- $R_s^a$  is the immediate reward after transitioning from state s to s' using action a.  
Where  $R_s^a = E[R_{t+1} | S_t = s, A_t = a]$
- $\gamma \in [0,1]$  is the discount factor applied to the reward

For a state to be Markov it needs to satisfy the following condition:

$$P[S_{t+1} | S_t] = P[S_{t+1} | S_1, \dots, S_t]$$

This means that the current state is required to contain all the information of the previous states. For a MDP all states must be Markov.

We now define the definition which is the total discounted reward from time-step t,

named the return  $G_t$  where:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots \quad (3.1)$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (3.2)$$

The discount  $\gamma \in [0,1]$  determines the present value of future rewards. For  $\gamma = 0$  the return  $G_t$  only depends on the current reward that can be obtained in the next step  $R_{t+1}$ . Which can be said to be "short-sighted". While for  $\gamma = 1$  the return depends on all the rewards that are projected to be obtained until the process terminates. This can be said to be "far-sighted". The larger  $\gamma$  is the more the rewards of later steps closer to the terminating state affects the return.

We now define another important required definition, namely the state value function  $v(s)$  where:

$$v(s) = E[G_t | S_t = s]$$

We can also decompose  $v(s)$  so that it becomes a recursive function as follows:

$$v(s) = E[G_t | S_t = s] \quad (3.3)$$

$$= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \quad (3.4)$$

$$= E[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) | S_t = s] \quad (3.5)$$

$$= E[R_{t+1} + \gamma G_{t+1} | S_t = s] \quad (3.6)$$

$$= E[R_{t+1} + \gamma v(S_{t+1}) | S_t = s] \quad (3.7)$$

## **Chapter 4**

### **Summary and Conclusion**

# Bibliography

- [1] P. U. F. of Computer Science, “8puzzle assignment.” [Online]. Available: <https://www.cs.princeton.edu/courses/archive/spring18/cos226/assignments/8puzzle/index.html>
- [2] W. E. Johnson, Wm. Woolsey; Story, “Notes on the ”15” puzzle,” *American Journal of Mathematics*, vol. 29, no. 2(4), p. 397–404, 1879.
- [3] D. of the Cube Forum, “5x5 sliding puzzle can be solved in 205 moves,” 2016. [Online]. Available: <http://cubezzz.dyndns.org/drupal/?q=node/view/559>
- [4] Minsky, *Computation: finite and infinite machines*. Prentice Hall, 1967.
- [5] R. S. Sutton and A. G. Barto, *Reinforcement Learning An Introduction second edition*. The MIT Press, 2018.

# **Appendix A**

## **Project Planning Schedule**

This is an appendix.



# **Appendix B**

## **Outcomes Compliance**



This is another appendix.

# Appendix C

## Student and Supervisor agreement

### Agreement between skripsie student and study leader regarding mutual responsibilities

Project (E) 448, Department of Electrical and Electronic Engineering, Stellenbosch University

Student name and SU#:	Umr Barends 18199313
Study leader:	Mr JC Schoeman
Project title:	Learning to solve Sliding Puzzles using Reinforcement Learning
Project aims:	Generally in robotics a manipulator (eg. an arm) is used to manipulate an object in the environment. It is normally easier to first simulate the robots behavior in a more simple environment. In this project we try to solve a sliding puzzle using reinforcement learning, where the same algorithm can then later be applied to the robotics problem.
	<ol style="list-style-type: none"><li>1. It is the responsibility of the student to clarify aspects such as the definition and scope of the project, the place of study, research methodology, reporting opportunities and -methods (e.g. progress reports, internal presentations and conferences) with the study leader.</li><li>2. It is the responsibility of the study leader to give regular guidance and feedback with regard to the literature, methodology and progress.</li><li>3. The rules regarding handing in and evaluation of the project is outlined in the Study Guide/website and will be strictly adhered to.</li><li>4. The project leader conveyed the departmental view on plagiarism to the student, and the student acknowledges the seriousness of such an offence.</li></ol>
Signature — study leader:	
Signature — student:	
Date:	5 August 2020

(Upload this form to SUNLearn)

Figure C.1: Student and Supervisor agreement