



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

November 1, 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 / November 1, 2020

Abstract

English

The English abstract.

Contents

Declaration	ii
Abstract	iii
List of Figures	vi
List of Tables	viii
Nomenclature	ix
1. Introduction	1
1.1. Problem Statement and Project Objective	1
1.2. Background	1
1.3. Document Outline	2
2. Literature Review	3
2.1. Introduction	3
2.2. Solving sliding puzzle using value iteration	3
3. Markov Decision Processes	4
3.1. Introduction	4
3.2. Definition of a MDP	4
3.3. Value function and policies	5
3.4. Bellman Expectation Equation	6
3.5. Optimal Policies and Bellman Optimality Equations	7
4. Dynamic Programming	8
4.1. Policy Evaluation	8
4.2. Policy Improvement	9
4.3. Policy iteration	10
4.4. Value Iteration	10
4.5. Value iteration grid-world example problem	11
5. Reinforcement Learning	13
5.1. Model-free prediction	13
5.1.1. Monte Carlo prediction	13

5.1.2. Temporal Difference (TD) Learning	14
5.2. Model-free control	15
5.2.1. Exploration and exploitation	16
5.2.2. SARSA: On policy TD control	16
5.2.3. Q-learning: Off policy TD control	18
6. Applying reinforcement learning to sliding puzzles	20
6.1. Introduction	20
6.2. SARSA figs	22
6.3. Q-learning figs	23
6.4. Discussion	23
7. Summary and Conclusion	24
7.1. Summary	24
7.2. Future work and recommendations	24
Bibliography	25
A. Solvability of a NxN sliding puzzle	26
A.1. General puzzle description	26
A.2. Conditions for solvability	26
B. Project Planning Schedule	28
C. Outcomes Compliance	29
D. Student and Supervisor agreement	30

List of Figures

1.1. Figures of a 3x3 sliding puzzle being solved	2
4.1. Policy iteration with images inspired by Sutton and Barto [1]	9
4.2. This is a grid world example solved using value iteration. The goal is to have values in the cells represent the amount of steps it would take from that cell to the cell in the top left corner of that iteration. Each step results in a reward of -1. All iterations after iteration $k = 7$ are the same, although not shown, as the optimal state-value function $v_7(s) = v_*(s)$ has been reached by then.	11
5.1. SARSA backup diagram representation [2]	17
5.2. Q-learning backup diagram where the red line indicates a max across the actions that can be taken from state S' [2]	18
6.1. Green tiles can have any value, as long as they don't have the same value as any other tile.	20
6.2. Illustrating where the blank tile is allowed to move for agent 2. The blank tile can move into the blue, green, number 2 and 3 tiles, but cannot move into any of the red tiles positions. For the configuration to be solved the blue tiles cannot contain a blank tile, nor a number 2 or 3 tile.	21
6.3. Illustrating where the blank tile is allowed to move for agent 3. The blank tile can move into the green, number 2 and 3 tiles, but cannot move into any of the red tiles positions. When the 2 and 3 numbered tiles are placed in the positions as seen in this figure, agent 4 reaches its terminal state. . .	21
6.4. Illustrating where the blank tile is allowed to move for agent 4. The blank tile can move into the green, number 4 and 7 tiles, but cannot move into any of the red tiles positions.	22
6.5. Illustrating where the blank tile is allowed to move for agent 4. The blank tile can move into the green, number 5, 6 and 8 tiles, but cannot move into any of the red tiles positions.	22
6.6. SARSA episode vs return graph for agent 1, using 100 runs to average. ϵ is the only hyper parameter varied, with $\gamma = 0.99$ and $\alpha = 0.1$ kept constant.	22
A.1. Example of a sliding puzzle	26
A.2. Odd boards with change in blank piece only having even inversion change [3]	27

A.3. Even board solvability [3]	27
D.1. Student and Supervisor agreement	30

List of Tables

Nomenclature

Variables and functions

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

Acronyms and abbreviations

RL	Reinforcement Learning
SARSA	State-Action-Reward-State=Action
Q-learning	Q value learning
MDP	Markov Decision Process

Chapter 1

Introduction

1.1. Problem Statement and Project Objective

In robotics, a manipulator (ie. a robotic arm) is a device used to manipulate objects in its environment. When the manipulation task is relatively complex, it is often difficult or even impossible for a human to direct how the robot should act. One possible solution is that the robot learn by itself which actions are best, which can be done through reinforcement learning. A popular method for doing developing reinforcement learning algorithms for robotics is to first solve similar, less complex problems.

In this project reinforcement learning is used to solve sliding puzzle. An example of a sliding puzzle being solved can be seen in Figure 1.1. The objective of this project is to solve a shuffled puzzle with the minimum amount of moves using reinforcement learning. The rules of the game are explained in the next section.

1.2. Background

The sliding puzzle is a game with a history dating back to the 1870's [4]. The game consists of $N^2 - 1$ tiles arranged on an $N \times N$ grid with one empty tile. The tiles are numbered 1 to N . A possible configuration of a 3×3 puzzle is shown in Figure 1.1a. The puzzles state can be changed by sliding one of the numbered tiles, next to the empty tile, into the place of the empty tile. The empty tile then takes the place of the numbered tile. This can visually be seen between Figures 1.1a and 1.1b where tile number 6 moves right in the transition between Figures 1.1a and 1.1b. We will denote the action that can occur in a certain state by the possible movements the empty tile can move. Which will be defined as $A_s = up, down, left, right$.

The sliding puzzle is a game with a history dating back to the 1870's [4]. Many solutions have been thought up to solve such puzzles, which include using search algorithms [5]. The rules of the game are that the only tiles that can move are the ones next to the tile with no number, known as the 'blank tile'. Additionally only one tile can be moved at a time. The final puzzle configuration is shown in 1.1d, where all the tiles are placed in ascending order read left to right and top to bottom. The aim of the game is to arrange a

shuffled puzzle into the final puzzle configuration by moving one tile at a time.

1	2	3
4	6	
7	5	8

(a) Puzzle with tiles 5,6,8 and blank in the incorrect positions

1	2	3
4		6
7	5	8

(b) Puzzle with tiles 5,8 and blank in the incorrect positions

1	2	3
4	5	6
7		8

(c) Puzzle with tiles 8 and blank in the incorrect positions

1	2	3
4	5	6
7	8	

(d) Solved puzzle

Figure 1.1: Figures of a 3x3 sliding puzzle being solved

1.3. Document Outline

In Chapter 2 we will look at a paper which used an alternate approach to solve a 16 tile sliding puzzle. Chapter 3 describes in detail theory relating to Markov decision processes. Chapter 5 discusses in detail the dynamic programming and reinforcement learning techniques which will be used to solve our puzzle problem.

In Chapter ?? the theory of reinforcement learning is applied to solve sliding puzzles. In the same chapter we

Rather how you used the theory of RL to design the solution to the sliding puzzle problem. Do this mathematically. Talk about things like the size of the state space, the sequential agents that you trained, how you decided on algorithms or hyperparameters.

Chapter 6 contains tests and verification that our solution works using the results of various experiments. Chapter 7 concludes the report.

Chapter 2

Literature Review

2.1. Introduction

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. [1] This path is typically described as a policy π .

learning and planning lecture slide 42 (37/46)

2.2. Solving sliding puzzle using value iteration

In this section a few terms will be used which are fully described in the theory chapters 3 and 5. These terms include value iteration and policy iteration.

A 4x4 puzzle solution is investigated in [6]. The method they use is a variation of value iteration. The amount of states in a 4x4 puzzle is approximately 10^{13} elements. It was found in [6] that value iteration is indeed a feasible approach if the value iteration algorithm is altered so that the only a subset of the entire state space is used sequentially until a solution is found.

This reducing of the state space is similar to what will later be done in this report. However the method we use to reduce the state space is by having our algorithm be guided somewhat by how a human would solve a puzzle, by doing the edges first. Additionally that we use a different reinforcement learning method known as policy iteration.

Solutions to sliding puzzles can also been found by using search algorithms such as A* and IDA* [5].

Chapter 3

Markov Decision Processes

3.1. Introduction

In order to describe reinforcement learning techniques, we first have to have a way of mathematically defining reinforcement learning problems. In this chapter, we do so by formally introduce the concept of Markov decision processes (MDPs). This chapter introduces the key concepts required to address the reinforcement learning problem. These concepts includes, state and action space, returns, value functions and Bellman equations. The majority of the theory discussed in this chapter is based on the work by Sutton and Barto [1].

3.2. Definition of a MDP

A Markov decision process is a stochastic, discrete, time controlled process. What this means is that a MDP is a partially controlled and partially random process [1]. In general RL problems can be described using MDPs [2]. Paul Weng states that, MDP's are described by the following: state-space S , action-space A , state transition probability $P_{ss'}^a$, reward function R_s^a and discount factor γ [7]. These variables can be denoted in a 5-tuple as

$$(S, A, P_{ss'}^a, R_s^a, \gamma) \quad (3.1)$$

where,

- S is a finite set of states with state $s \in S$
- A is a finite set of actions with action $a \in A$
- $P_{ss'}^a$ is a matrix of probabilities which predicts the next state s' , where $P_{ss'}^a = P(S_{t+1} = s' | S_t = s, A_t = a)$
- 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

In an MDP, the assumption is that all the states have the Markov property [1]. This means that

$$P[S_{t+1}|S_t] = P[S_{t+1}|S_1, \dots, S_t]. \quad (3.2)$$

In other words, the probability of transitioning to state S_{t+1} given state S_t , must be equal to the probability of transitioning to state S_{t+1} given all states S_1 to S_t . This means that the current state is required to contain all the information of the previous states.

3.3. Value function and policies

The return G_t is the total sum of the reward at each time step, multiplied by a constant diminishing factor γ . The discount $\gamma \in [0,1]$ determines the present value of future reward R_t .

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \end{aligned} \quad (3.3)$$

For $\gamma = 0$, the return G_t only depends on the single reward that can be obtained in the next step, namely R_{t+1} . When $\gamma = 0$ the return G_t is 'short-sighted'. For $\gamma = 1$ on the other hand, the return depends on all the rewards that are projected to be obtained until the process terminates. G_t is thus "far-sighted" for $\gamma = 1$ [1]. What can also note is that using a discount factor between 0 and 1 makes the return finite because R_{t+k+1} is finite and

$$\sum_{k=0}^{\infty} \gamma^k = \frac{1}{1 - \gamma}, \quad (3.4)$$

which is finite for $\gamma \in [0,1]$.

The probability that the agent will choose a specific action given the state is called its policy $\pi(a|s)$. Mathematically this is defined as

$$\pi(a|s) = P[A_t = a, S_t = s]. \quad (3.5)$$

Another important definition is the value function $v(s)$ defined as

$$v(s) = E[G_t | S_t = s]. \quad (3.6)$$

This equation shows that the value function $v(s)$ is the expected return given that $v(s)$ is evaluated in state s at time-step t .

3.4. Bellman Expectation Equation

Sutton and Barto [1] decomposes $v(s)$ so that it becomes a recursive function as follows

$$v(s) = E[G_t | S_t = s] \quad (3.7)$$

$$= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \quad (3.8)$$

$$= E[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) | S_t = s] \quad (3.9)$$

$$= E[R_{t+1} + \gamma G_{t+1} | S_t = s] \quad (3.10)$$

$$= E[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]. \quad (3.11)$$

This means that $v(s)$ only depends on R_{t+1} (the reward that can be obtained in the next step) and $\gamma v(S_{t+1})$, the discounted value function at time $t+1$. What has been defined in equation 3.11 is known as the Bellman equation. [1] Sutton and Barto then defines the value function given that the agent follows the policy π as

$$v_\pi(s) = E_\pi[G_t | S_t = s]. \quad (3.12)$$

Additionally, they define the action-value function (q-function),

$$q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a]. \quad (3.13)$$

The q-function is defined as the expected return of taking action a from state s and thereafter following policy π . What can be noted is that $v(s)$ is a prediction of what the value for being in a certain state is, which is related to planning. While $q(s, a)$ is the value for being in a certain state and taking an action, which relates to control.

Once again, Sutton and Barto [1] decomposes the state-value and action-value functions to be in the form of the Bellman Equation in equation 3.11. This results in the Bellman *expectation* equations

$$v_\pi(s) = E_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s] \quad (3.14)$$

$$= \sum_{a' \in A} \pi(a' | s) (R_s^{a'} + \gamma \sum_{s' \in S} P_{ss'}^{a'} v_\pi(s')) \quad (3.15)$$

and

$$q_\pi(s, a) = E_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a] \quad (3.16)$$

$$= R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \sum_{a' \in A} \pi(a' | s') q_\pi(s', a'). \quad (3.17)$$

Equation 3.15 can be placed into vector form as:

$$v_\pi = R^\pi + \gamma P^\pi v_\pi \quad (3.18)$$

Which has the solution:

$$v_\pi = (I - \gamma P^\pi)^{-1} R^\pi \quad (3.19)$$

Equation 3.19 is a concise formula which is easily implemented in code.

3.5. Optimal Policies and Bellman Optimality Equations

Sutton and Barto now defines the optimal policy, which only has value 1 if the action the agent takes maximizes $q_*(s, a)$ [1]. This is mathematically defined as:

$$\pi_*(a|s) = \begin{cases} 1, & \text{if } a = \arg \max_{a \in A} q_*(s, a) \\ 0, & \text{otherwise} \end{cases} \quad (3.20)$$

There always exists an optimal deterministic policy for any MDP [1]. We only need to know $q_*(s, a)$ to know the optimal policy. This is important because the optimal policy describes which actions must be taken to maximize rewards, which is precisely the goal of RL. Sutton and Barto [1] then defines what is know as the Bellman *optimality* equations.

$$v_*(s) = \max_a (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s')) \quad (3.21)$$

and

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_{a'} q_*(s', a'). \quad (3.22)$$

What equation 3.21 represents is the maximum return that can be obtained from state s , when following the optimal policy π_* . Equation 3.22 represents the maximum return that can be obtained by taking action a at state s and thereafter following the policy. At this point, we have outlined a mathematical background for reinforcement learning problems. In the next two chapters we will discuss methods of solving the MDP and obtaining the optimal state-value and action-value functions.

Chapter 4

Dynamic Programming

Dynamic programming (DP) is a method to solve problems by using a divide and conquer approach [2]. That is to say a problem is broken into smaller sub-problems and then solved. Dynamic programming is used to solve MPDs of which we have perfect models. Most times, we do not have access to such models and hence DP is of limited use in reinforcement learning. However, DP is still of theoretical importance and a useful component to understanding reinforcement learning [1]. There are two properties that must be met for a problem to be solvable via DP, namely optimal substructure and overlapping sub-problems [2]. **Optimal substructure** requires that the solution can be decomposed into sub-problems. **Overlapping sub-problems** requires that the sub-problems recur, so that the solution can be cached and reused.

Markov decision processes satisfy both of these conditions. According to David Silver, the Bellman equation provides a recursive decomposition which fulfills the optimal substructure property, while the value function caches and reuses solutions satisfying the overlapping sub-problem's condition [2]. The key concept of DP, is to use state-value functions as the structure through which optimal policies are found [1]. As a basis to finding the optimal policy, the next section first discusses how to evaluate a given policy.

4.1. Policy Evaluation

In order to evaluate a given policy, we calculate the state-value function V_π for a given policy π . Sutton and Barto names this policy evaluation [1]. From Equation 3.15 we know that,

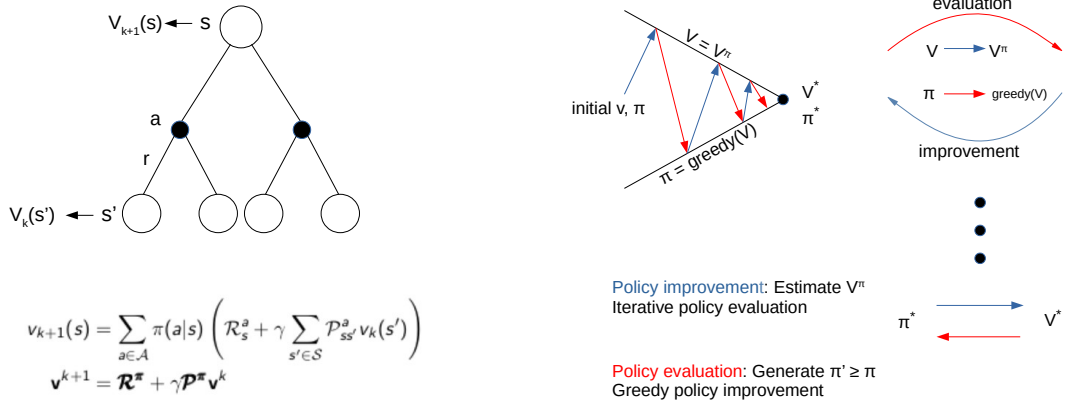
$$v_\pi(s) = \sum_{a' \in A} \pi(a'|s) (R_s^{a'} + \gamma \sum_{s' \in S} P_{ss'}^{a'} v_\pi(s')).$$

continue editing from here In order to make it clear how policy evaluation works, we use a diagram inspired by David Silvers lecture series [2] seen in Figure 4.1a. Starting at the root node, state s in the diagram in Figure 4.1a, the agent can take any action based on its policy π . After the agent takes this action, the environment responds and moves the agent to state s' . The environments response depends on the probabilities contained in the function $P_{ss'}^a$.

Sutton and Barto then introduces the notation for an approximate value function, $v_k(s)$ which maps every state $s \in \mathcal{S}$ to a real number in \mathbb{R} [1]. They then choose v_0 arbitrarily, where after $v_k(s)$ is successively updated using

$$v_{k+1}(s) = \sum_{a' \in A} \pi(a'|s) (R_s^{a'} + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^{a'} v_k(s')) \quad (4.1)$$

The value function is updated continuously in this way, until the difference between v_k and v_{k+1} is determined negligible. This process is known as iterative policy evaluation. Now that we have a tool to evaluate a policy, in the next section we discuss how to improve a given policy, in order to find the optimal policy π_* .



4.2. Policy Improvement

According to David Silver, a given policy π can be improved using the following two steps [2]:

Step one is *evaluating* the policy π using equation 3.15:

$$v_\pi(s) = E_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$$

Step two is *improving* the policy by acting greedy with respect to v_π to obtain a new policy π' .

Sutton and Barto defines acting greedy [1], as the greedy policy which makes the agent

always takes the action that looks best according to v_π looking one step ahead,

$$\pi' = \text{greedy}(v_\pi). \quad (4.2)$$

Sutton and Barto then further expand the greedy policy concept to consider all states and actions resulting in

$$\pi'(s) = \arg \max_a (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s')), \quad (4.3)$$

where $\arg \max_a$ denotes the action which if taken maximizes the expression following it. By using the greedy policy, Sutton and Barto proves in [1] that we will always obtain the optimal expected return v_* defined in equation 3.21.

Sutton and Barto provides a proof that using this method always results in a policy π' better, or as good as π , which will not be provided here but can be found in [1]. That is to say the expected return following π' will be greater or equal than following π

$$v_{\pi'} \geq v_\pi. \quad (4.4)$$

In the next two sections we will use policy improvement as a tool to find the optimal policy via two methods, policy iteration and value iteration.

4.3. Policy iteration

There are two separate uses for dynamic programming in reinforcement learning. These are prediction and control which both have (S, A, P, R, γ) as input.

- The prediction method uses (S, A, P, R, γ) and a policy π as input, while it outputs a value function v_π . This is known as value iteration.
- The control method only takes (S, A, P, R, γ) as input and outputs the optimal value function v_* and optimal policy π_* . This is known as policy iteration.

Figure 4.1b illustrates how greedy policy iteration works. The blue lines represents policy improvement and the red lines policy evaluation. By continuously using these two steps in succession, we eventually converge to the optimal state-value function v_* and optimal policy π_* [1]. It is known that policy iteration converges, as it uses policy improvement which we know converges from the previous section [1].

4.4. Value Iteration

Policy iteration has a drawback of needing for every iteration, to do policy evaluation which in itself is a iterative computation requiring many calculations across the entire

state-space [1]. Iterative policy evaluation only converges in the limit of its time steps. However Sutton and Barto questions if it is necessary to converge precisely on the optimal policy, or if we can stop reasonably short and acquire an approximate optimal policy [1]. They answer this by introducing the concept of value iteration, where they use policy iteration with a special case of policy evaluation as follows:

$$v_{k+1}(s) = \max_{a \in A} (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_k(s')) \quad (4.5)$$

and in vector form,

$$\mathbf{v}_{k+1} = \max_a (\mathbf{R}^a + \gamma \mathbf{P}^a \mathbf{v}_k) \quad (4.6)$$

Note these equations are the Bellman optimality equation (in equation 3.21), turned into an update rule. In the next chapter we use almost the same method to find the optimal policy, except that we do not assume that we have access to a perfect model of the MDP.

4.5. Value iteration grid-world example problem

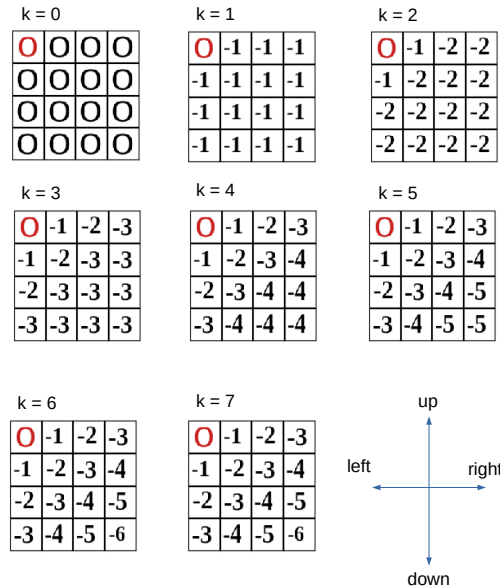


Figure 4.2: This is a grid world example solved using value iteration. The goal is to have values in the cells represent the amount of steps it would take from that cell to the cell in the top left corner of that iteration. Each step results in a reward of -1. All iterations after iteration $k = 7$ are the same, although not shown, as the optimal state-value function $v_7(s) = v_*(s)$ has been reached by then.

Figure 4.2 is an example problem, which David Silver [2] uses to illustrate how value

iteration works. The states are defined as cells contained within a grid. The example in Figure 4.2 has a state space of 16 states with $S \in [0,15]$ and an action space of actions $A = \{\text{up, down, left, right, NONE}\}$ indicated on the figure. The action NONE is used to indicate the agent choosing to stay in its current state. Each state consists of an index and a value, with the terminal state being at index 0. The agent can move in all directions in all states, except for the terminal state, where it can only choose to stay still. The state indices are numbered from left to right and top to bottom in ascending order starting at 0. Hence we can treat the states as an array in computer programming with 0 to 3 being the top row indices, 3 to 7 the second most top row indices etc.

Imagine that we have an agent dropped into a grid as is shown in the initial iteration in Figure 4.2 then

$$v_0(0) = v_0(1) = v_0(2) = v_0(3) = v_0(4) = \dots = v_0(s) = 0. \quad (4.7)$$

Hence at iteration 0 no matter which state the agent is dropped into, all actions are equally valuable to it. As the iterations k increases, $v_k(s) \rightarrow v_*(s)$. We know the agent has learnt the optimal state-value function, when the current iterations state-value function has a negligible, or no difference from its previous iteration

$$v_{k+1}(s) - v_k(s) \leq \delta, \quad (4.8)$$

where δ is a chosen threshold.

Chapter 5

Reinforcement Learning

In previous chapters we discussed methods which required knowledge about the model of the environment. This knowledge was embedded in the probability transition matrix P . Often times it is not possible to know the model of the environment beforehand. This is where reinforcement learning is useful because it does not require a model of the environment.

In reinforcement learning, the way an agent learns the best actions to take, is by experiencing sequences of states and actions, defined as an episode. When the agent reaches the final state, we say the episode has terminated. In this chapter we discuss two methods of reinforcement learning (Monte Carlo and temporal difference methods), both of which can be implemented via prediction or control methods. We will not cover the Monte Carlo control method, but it can be found in Sutton and Barto [1]. We begin in the next section with the prediction methods.

5.1. Model-free prediction

5.1.1. Monte Carlo prediction

According to Sutton and Barto the term Monte Carlo is generally used for any estimation technique where a large portion of its operation is random [1]. Sutton and Barto specifically describes Monte Carlo methods as solving reinforcement learning problems by averaging sample returns [1]. Monte-Carlo methods require only experience, which it obtains by sampling sequences of states, actions and rewards. By sampling enough times, by the law of large numbers, eventually an accurate state-value function will be obtained [1].

Monte-Carlo policy evaluation works as follows. We define a total return S in state s as $S(s)$ and a integer counter $N(s)$ which represents the amount of times s is seen. We then initialize both $S(s)$ and $N(s)$ initialized to zero. Then at time step t and state s we increment the counter $N(s) = N(s) + 1$. This can be done in two ways, one being only incrementing $N(s)$ the *first* time state s is seen, the other being incrementing $N(s)$ *every* time states s is seen [1].

In this chapter we use the method of incrementing the counter every instance that state s is seen. Then we increment the total reward $S(s) = S(s) + G_t$ every time the agent

reaches state s . Finally we calculate the state value estimate function as $V(s) = \frac{S(s)}{N(s)}$. Then $V(s) \rightarrow V_\pi(s)$ as $N(s) \rightarrow \infty$, $\forall s \in S$. In other words when we have seen all states infinitely many times, we will have obtained exactly $V_\pi(s)$. Although theoretically we need to visit each state infinitely many times, practically there is a finite time-step where $V(s)$ converges to $V_\pi(s)$. In order to describe incremental Monte-Carlo updates, David Silver introduces the concept of an incremental mean u_k [2]. The derivation for u_k from a standard mean is as follows:

$$u_k = \frac{1}{k} \sum_{j=1}^k x_j \quad (5.1)$$

$$\begin{aligned} &= \frac{1}{k} \left(x_k + \sum_{j=1}^{k-1} x_j \right) \\ &= \frac{1}{k} (x_k + (k-1)u_{k-1}) \\ &= u_{k-1} + \frac{1}{k} (x_k - u_{k-1}). \end{aligned} \quad (5.2)$$

Since $V(s)$ is a mean, we break it up similarly to Equation 5.2, with $u_{k-1} = V(s)$, $N(s) = t = k$ and return $G_t = x_k$. Substituting these variables into equation 5.2 yields that at every time-step t :

$$N(S_t) \leftarrow N(S_t) + 1 \quad (5.3)$$

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t)). \quad (5.4)$$

Sutton and Barto in [1] sets $\frac{1}{N(S_t)} = \alpha$, where $0 < \alpha < 1$ resulting in

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t)). \quad (5.5)$$

Equation 5.5 gives us a way to find the optimal state-value function $V_\pi(S_t)$ by iteratively updating $V(S_t)$. Sutton and Barto names Equation 5.5, *constant* $-\alpha$ Monte Carlo.

5.1.2. Temporal Difference (TD) Learning

Temporal difference (TD) learning uses a combination of Monte Carlo and dynamic programming (DP) ideas. TD is one of the unique aspects of reinforcement learning, unlike Monte Carlo and DP which have more general applications besides reinforcement learning [1]. Like Monte Carlo TD learns without requiring a model of the environment. Because TD is a reinforcement method, it can be split into prediction and control methods as was stated in the introduction of this chapter. Here in this section we only describe the prediction method, while the next section will discuss two TD control methods, namely

SARSA and Q-learning. The aim of TD learning is to learn a value function v_π from experience by following policy π . The simplest form of TD is replacing the return G_t in Equation 5.5 from the previous section as follows:

$$V(S_t) = V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)). \quad (5.6)$$

Where $[R_{t+1} + \gamma V(S_{t+1})]$ is known as the TD target, the value which $V(S_t)$ tends to as $t \rightarrow \infty$ and $[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$ is known as the TD error. The aim of TD learning is to reduce the TD error to zero and update $V(S_t)$ to the estimated return $R_{t+1} + \gamma V(S_{t+1})$. Sutton and Barto mathematically describes the relationship between MC and TD using what they define as the n-step return $G_t^{(n)}$ described by:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n}). \quad (5.7)$$

For clarity we show $G_t^{(n)}$ for a few values of n as follows:

$$(TD)n = 1 : G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1}) \quad (5.8)$$

$$n = 2 : G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2}) \quad (5.9)$$

$$(MC)n = \infty : G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T. \quad (5.10)$$

What can be seen is, as $n \rightarrow \infty$, the return in TD learning turns into the same form as in Monte Carlo. Hence Monte Carlo can be said to be a special case of TD learning. Substituting $G_t^{(n)}$ into Equation 5.6 yields what Sutton and Barto calls n-step TD as:

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t^{(n)} - V(S_t)). \quad (5.11)$$

What we now have is a method to control how many steps we want an agent to look ahead. This is useful because $V(S_t)$ can converge more quickly to $V_\pi(S_t)$ if n is chosen well. Normally one has to test for different values of n to see which value works best for the specific problem at hand. In the next section we no longer look at estimating the state-value function $V_\pi(S_t)$, but instead work with *state-action* pairs in order to find the optimal state-action function $q_*(s, a)$.

5.2. Model-free control

Model-free control can be implemented in two ways. **On-policy learning**, which lets the agent learn characteristics of a policy π by sampling from π and **off-policy learning**, learning characteristics of policy π by sampling another policy.

In order to find the optimal state-action function $q_*(s, a)$, we require a method of evaluating

a state-action policy. Previously, when using greedy policy improvement over the state-value function $V(s)$ we had in Equation 4.2:

$$\pi'(s) = \max_{a \in A} (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s'))$$

Sutton and Barto applies greedy policy improvement over the state action function $Q(s,a)$ as follows:

$$\pi'(s) = \max_{a \in A} Q(s, a). \quad (5.12)$$

One should notice that to find the optimal policy (by using Equation 5.12 for policy improvement), we no longer need a model of the environment. In other words there is no dependency on the probability matrix $P_{ss'}^a$. In the next section we will see that there is a problem if an agent only uses Equation 5.12 for policy improvement. In the same section a solution is then given to the problem.

5.2.1. Exploration and exploitation

In reinforcement learning, exploration refers to an agent taking a random action and exploitation when the agent acts greedily. When an agent exploits it can obtain a larger return and when it explores it gets a more accurate estimate of the action-value function. The agent can not simultaneously explore and exploit which results in what is known as the exploration-exploitation dilemma [8]. To ensure that the agent explores occasionally instead of always taking the greedy action with respect to the state-action function $Q(s,a)$, we let the agent act randomly with a probability ϵ . This can mathematically be expressed as:

$$\pi(a|s) = \begin{cases} 1 - \epsilon, & \text{if } a^* = \arg \max_{a \in A} Q(s, a) \\ \epsilon, & \text{otherwise take a random action} \end{cases} \quad (5.13)$$

Equation 5.13 is known as the ϵ -greedy policy. Using the ϵ -greedy policy results in the agent exploring occasionally, which allows for faster convergence if ϵ is chosen appropriately. Usually, an ϵ of 0.1 is a good starting value.

5.2.2. SARSA: On policy TD control

We now look at using TD methods for control. The backup diagram for SARSA is shown in Figure 5.1, where we visually see a sequence of a State-Action-Reward-State-Action (SARSA). This is where the name SARSA originates from. By following the same process that we used in the TD prediction section, we can arrive to an equation similar to Equation 5.6. Instead of updating the state-value function $V(S)$, we update the state-action function

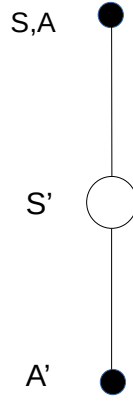


Figure 5.1: SARSA backup diagram representation [2]

$Q(S, A)$. For every state action pair, we can calculate the state-action function $Q(S, A)$ using TD learning as follows:

$$Q(S, A) = Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]. \quad (5.14)$$

Where $R + \gamma Q(S', A')$ is the TD target. Similar to Monte Carlo prediction where a n-step state value return G_t^n was defined in Equation 5.7, Sutton and Barto also defines n-step Q-return as:

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n}). \quad (5.15)$$

Although Sutton and Barto do not explicitly mention it, but $q_t^{(n)}$ is assumed to be zero for all $n < 1$. The goal of n-step SARSA is to update $Q(s, a)$ towards the n-step Q-return $q_t^{(n)}$,

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha[q_t^{(n)} - Q(S_t, A_t)]. \quad (5.16)$$

Sutton and Barto states that SARSA always converges towards an optimal action-value function, as long as all state-action pairs are visited enough times [1]. One can also deduce this fact from the law of large numbers, as was stated in the Monte Carlo prediction section.

Algorithm 1 below, adapted from Sutton and Barto, shows the pseudo-code which can be used to implement SARSA. Lines 4 and 7 in algorithm 1 is policy improvement. Line 8 is the SARSA update line corresponding to Equation 5.16, with $n=1$. Later in chapter 6 we will apply this algorithm, where will adjust the hyper parameters to see how they affect the agents learning process.

Algorithm 1: SARSA algorithm to find optimal state action function $Q_*(S, A)$

function SARSA (α, γ, ϵ);

Input : Three nonnegative integers, learning rate α , discount γ and exploration probability ϵ

Output: $Q(S, A)$

- 1: Initialize $Q(s,a)$, $\forall s \in S, a \in A(s)$ arbitrarily and $Q(\text{terminal state}, \{\text{all actions}\}) = 0$
 - 2: for each episode:
 - 3: Initialize S
 - 4: Choose A from S using ϵ -greedy policy derived from Q
 - 5: for each step in episode:
 - 6: Take action A , observe R, S'
 - 7: Choose A' from S' using ϵ -greedy policy derived from Q
 - 8: $Q(S,A) \leftarrow Q(S,A) + \alpha[R + \gamma(Q(S',A') - Q(S,A))]$
 - 9: $S \leftarrow S'; A \leftarrow A';$
 - 10: Until S is terminal
-

5.2.3. Q-learning: Off policy TD control

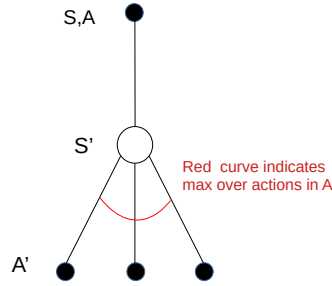


Figure 5.2: Q-learning backup diagram where the red line indicates a max across the actions that can be taken from state S' [2]

$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma \max_{a'} Q(S', a') - Q(S, A)). \quad (5.17)$$

We describe now Q-learning which falls under off policy learning. What we have been dealing with up to this chapter, was on-policy learning. Q-learning is an off-policy method, because it not only updates $Q(S,A)$, but it also updates the TD target $[R + \gamma \max_{a'} Q(S', a')]$. By comparing Figure 5.2 and Figure 5.1, we can see that the only difference between Q-learning and SARSA is that Q-learning takes the action a' , which results in the maximum state-action function $Q(S', a')$. We will also apply the Q-learning algorithm in Chapter 6. In this chapter we have covered various methods of reinforcement learning, of which we will only be applying TD control methods, namely SARSA and Q-learning. In the next chapter we will also discuss what effects each hyper parameter has on the agents learning process.

Algorithm 2: Q-learning algorithm to find optimal state action function $Q_*(S, A)$

function Q-learning (α, γ, ϵ);

Input : Three nonnegative integers, learning rate α , discount γ and exploration probability ϵ

Output : $Q(S, A)$

- 1: Initialize $Q(s,a)$, $\forall s \in S, a \in A(s)$ arbitrarily and $Q(\text{terminal state}, \{\text{all actions}\}) = 0$
 - 2: for each episode:
 - 3: Initialize S
 - 4: for each state in episode:
 - 5: Choose A from S using policy derived from Q (eg., ϵ -greedy)
 - 6: Take action A , observe R, S'
 - 7: $Q(S,A) \leftarrow Q(S,A) + \alpha[R + \gamma \max_{a \in A} (Q(S',a) - Q(S,A))]$
 - 8: $S \leftarrow S'$;
 - 9: Until S is terminal
-

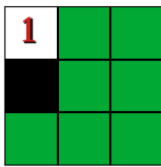
Chapter 6

Applying reinforcement learning to sliding puzzles

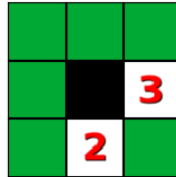
6.1. Introduction

In the following section we will test the SARSA and Q-learning algorithms described in chapter 5. The way the tests will be conducted is by varying the hyper parameters. These hyper parameters include, the discount factor γ , epsilon the probability of choosing a random action and α the learning rate.

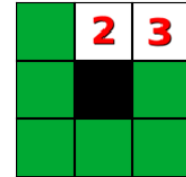
The puzzle that we will solve is of size 3x3. The way in which it will be solved is by training separate agents to solve different parts of the puzzle. The figures in Figure 6.1



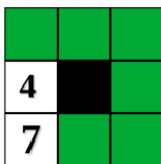
(a) First configuration that agent is trained to reach.



(b) Second configuration that agent is trained to reach. As long as the first tiles in the second and third row is not the blank, 2 or 3 numbered tiles, then the second configuration is solved.



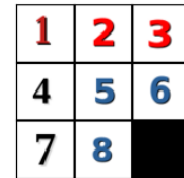
(c) Third configuration is solved when 2 and 3 numbered tiles are in the first row as shown.



(d) Fourth configuration is solved when 4 and 7 numbered tiles are in first column as shown.



(e) Fifth configuration is solved when the 2x2 block in the bottom-right of the grid is in order as shown.



(f) Fully solved puzzle.

Figure 6.1: Green tiles can have any value, as long as they don't have the same value as any other tile.

shows terminal states for each agent, in the order that we will use them, to reach the final

state in Figure 6.1f. For ease of referring to the different agents, we name the agent that solves the puzzle of the state in Figure 6.1a agent1, the agent of Figure 6.1b as agent2, the agent of Figure 6.1c as agent3, the agent of Figure 6.1d as agent4 and the agent of Figure 6.1e as agent5.

The amount of states that we have to look at when solving a puzzle, without breaking it up into different parts is $9! = 362880$ states. The optimal policy would be a certain episode, which is a sequence of states, actions and rewards. It would have to be seen multiple times in order for the agent to learn that, that policy is the best one.

When looking at the first configuration that the agent will be trained to solve in Figure 6.1a, there are only three unique tiles (blank, number 1 tile and green tiles) and there are 9 possible positions for a single tile. The blank is allowed to slide anywhere it wants according to the sliding puzzle games rules, state in the introduction chapter of this report. This results in a total state number of $9*8*7 = 504$ states.

For agents 2,3,4, and 5, the blank is no longer allowed to move anywhere on the grid, it is instead restricted to a $2*3$ or $3*2$ grid of movements. Agent 2 moves the blank, number 2, and number 3 tiles out of the first column if they are there. It does so by only looking at the second and third row as a $2*3$ puzzle to be solved. In other words the blank is only allowed to move in the last two rows as described in the caption of Figure 6.2. The state-space of agent 2 is $6*5*4 = 120$ states. Agent 3 works in a $2*3$ grid, contained by



Figure 6.2: Illustrating where the blank tile is allowed to move for agent 2. The blank tile can move into the blue, green, number 2 and 3 tiles, but cannot move into any of the red tiles positions. For the configuration to be solved the blue tiles cannot contain a blank tile, nor a number 2 or 3 tile.

the last two columns. The terminal state is shown in Figure 6.3. The terminal state and possible positions the blank tile can move of agent 4 is shown in Figure 6.4. The number of states of Agent 3 and 4 is the same as Agent 2 which is 120 states.

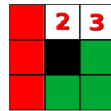


Figure 6.3: Illustrating where the blank tile is allowed to move for agent 3. The blank tile can move into the green, number 2 and 3 tiles, but cannot move into any of the red tiles positions. When the 2 and 3 numbered tiles are placed in the positions as seen in this figure, agent 4 reaches its terminal state.

Agent 5 is only allowed to move the blank tile in the bottom right 4 tiles as shown in Figure 6.5. There are only $4! = 24$ possible states that agent 5 can be in.

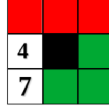


Figure 6.4: Illustrating where the blank tile is allowed to move for agent 4. The blank tile can move into the green, number 4 and 7 tiles, but cannot move into any of the red tiles positions.



Figure 6.5: Illustrating where the blank tile is allowed to move for agent 4. The blank tile can move into the green, number 5, 6 and 8 tiles, but cannot move into any of the red tiles positions.

If we use all agents (1 to 5 in ascending order), then the total amount of states that we have is $504 + (3 \cdot 120) + 24 = 888$ states. This is much smaller than if we had used brute force reinforcement learning, which would require us to deal with $9! = 362880$ states for the 3x3 sliding puzzle. We will now look at training results of SARSA being applied on the agents.

6.2. SARSA figs

Figure 6.6 depicts the results of using SARSA to learn to solve agent 1's problem in the previous section.

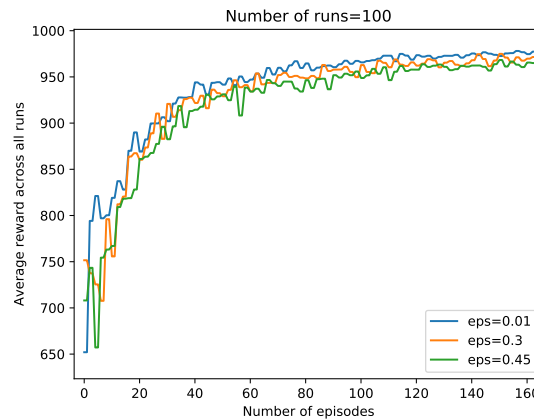


Figure 6.6: SARSA episode vs return graph for agent 1, using 100 runs to average. ϵ is the only hyper parameter varied, with $\gamma = 0.99$ and $\alpha = 0.1$ kept constant.

6.3. Q-learning figs

6.4. Discussion

Chapter 7

Summary and Conclusion

7.1. Summary

Objective 1 ...

Objective 2 ...

Objective 3 ...

So what? Relates back to problem and background sections in introductions.

7.2. Future work and recommendations

Bibliography

- [1] R. S. Sutton and A. G. Barto, *Reinforcement Learning An Introduction second edition*. The MIT Press, 2018.
- [2] D. Silver, “Ucl course on rl,” 2015. [Online]. Available: <https://www.davidsilver.uk/teaching/>
- [3] P. U. F. of Computer Science, “8puzzle assignment.” [Online]. Available: <https://www.cs.princeton.edu/courses/archive/spring18/cos226/assignments/8puzzle/index.html>
- [4] A. F. Archer, “A modern treatment of the 15 puzzle.” pp. 793–799, (1999).
- [5] R. E. Korf, “Depth-first iterative-deepening: an optimal admissible tree search. artificial intelligence 27,” pp. 97–109, (1985).
- [6] H. M. Bastian Bischoff, Duy Nguyen-Tuong and A. Knoll, “Solving the 15-puzzle game using local value-iteration,” 2013.
- [7] P. Weng, “Markov decision processes with ordinal rewards: Reference point-based preferences,” 2011.
- [8] M. Tokic, “Adaptive ϵ -greedy exploration in reinforcementlearning based on value differences,” 2014.
- [9] W. E. Johnson, Wm. Woolsey; Story, “Notes on the ”15” puzzle,” *American Journal of Mathematics*, vol. 29, no. 2(4), p. 397–404, 1879.
- [10] 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>
- [11] Minsky, *Computation: finite and infinite machines*. Prentice Hall, 1967.

Appendix A

Solvability of a NxN sliding puzzle

info in this chapter to still be integrated somewhere else into report if needed. Else will leave it out

A.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 A.1 we have a 1×16 array as: 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 A.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 A.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 A.1: Example of a sliding puzzle

A.2. Conditions for solvability

Even and odd sized boards are analysed separately (where 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. [3]

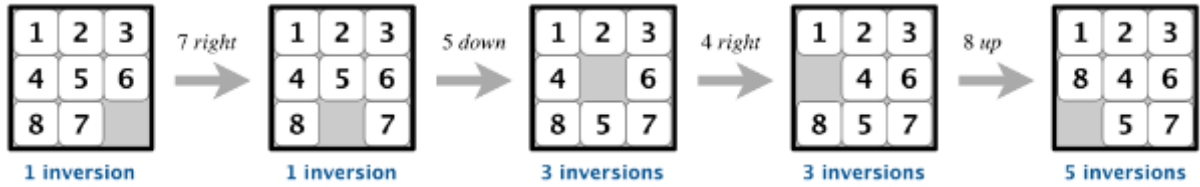


Figure A.2: Odd boards with change in blank piece only having even inversion change [3]

For even sized boards where N is even we have the board solvable if and only if the number of inversions plus the row of the blank square is odd. This is illustrated in Figure 3.

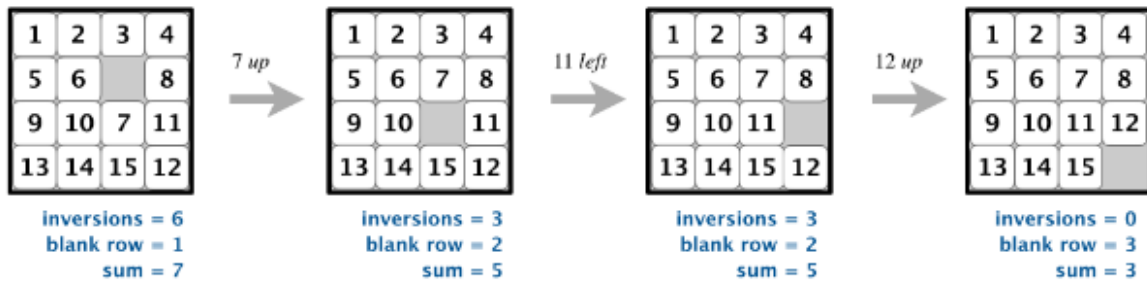


Figure A.3: Even board solvability [3]

Half of all puzzle configurations are unsolvable. [9] 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 [9]. Sliding puzzles can be solved relatively quickly with today's processing of computers for puzzles for example an 5×5 puzzle was solved in 205 tile moves in 2016. [10]

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. [11]

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.

Appendix B

Project Planning Schedule

This is an appendix.

Appendix C

Outcomes Compliance



This is another appendix.

Appendix D

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">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.2. It is the responsibility of the study leader to give regular guidance and feedback with regard to the literature, methodology and progress.3. The rules regarding handing in and evaluation of the project is outlined in the Study Guide/website and will be strictly adhered to.4. The project leader conveyed the departmental view on plagiarism to the student, and the student acknowledges the seriousness of such an offence.	
Signature — study leader:	
Signature — student:	
Date:	5 August 2020

(Upload this form to SUNLearn)

Figure D.1: Student and Supervisor agreement