



# Sorting Permutations by Reversals and Transpositions with Reinforcement Learning

*Guilherme Bueno Andrade      Andre Rodrigues Oliveira*  
*Zanoni Dias*

Relatório Técnico - IC-PFG-18-13

Projeto Final de Graduação

2018 - Julho

UNIVERSIDADE ESTADUAL DE CAMPINAS  
INSTITUTO DE COMPUTAÇÃO

The contents of this report are the sole responsibility of the authors.  
O conteúdo deste relatório é de única responsabilidade dos autores.

# Sorting Permutations by Reversals and Transpositions with Reinforcement Learning

Guilherme Bueno Andrade\*    Andre Rodrigues Oliveira†    Zanoni Dias†

## Abstract

*TODO*

## 1 Introduction

*TODO*

## 2 Reinforcement Learning

### 2.1 Overview

Reinforcement Learning (RL), like other branches of Machine Learning, has been drawing a lot of attention from the community in the recent years. Google DeepMind's AlphaGo victory over Lee Sedol [1] - world champion of the game Go, is one of many examples of recent astonishing applications of the technique. It consists of an agent learning how to accomplish a certain goal based on interactions with the environment.

Initially, the agent receives a state  $S_0$ . Based on that, the agent then takes an action  $A_0$ , ending up at state  $S_1$  and receiving some reward  $R_1$ . This process keeps going until the agent reaches a terminal state. Its goal is to maximize the total reward it gets along the way; i.e.,  $\max \sum_t R_t$ .

In order for the agent to accomplish such task, in value-based RL - the one being considered in this work, it optimizes the value function  $V(s)$ ,

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

where  $\gamma \in [0, 1)$  is a discount rate that makes the agent care more about most likely short term reward, and less about less probable future rewards.

---

\*gbuenoandrade@gmail.com

†Institute of Computing, University of Campinas, Brazil.

## 2.2 Exploitation vs. Exploration

A major concern in RL is the exploitation/explorarion trade-off. Exploration is about exploring new possibilities within the environment and finding out more information about it. Exploitation, on the other hand, is related to exploiting already known information so as to maximize the total reward.

Initially, the agent has no other option but to randomly explore the environment; however, as it learns about its surroundings, it can fall into the trap of sticking to safe known actions and miss larger rewards that depend on exploring unknown states.

This work uses the Epsilon-greedy strategy to address that problem. It specifies an exploration rate  $\epsilon$ , which is set to 1 initially. This rate defines the ratio of the steps that will be done randomly. Before selecting an action, the agent generates a random number  $x$ . If  $x > \epsilon$ , then it will select the best known action (exploit); otherwise, it will select an action at random (explore). As the agent acquires more knowledge about the environment,  $\epsilon$  is progressively reduced.

## 2.3 Q-table and the Bellman Equation

From the value function defined in equation 1, we have the Q-table, defined as follows.

$$Q(s, a) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s, A_t = a] \quad (2)$$

Where  $s'$  is the state that will be reached after the agent performs action  $a$ .

This is convenient because it allows the agent to pick the best action that can be performed from state  $s$  by simply finding  $\arg \max_a Q(s, a)$ .

Furthermore,  $Q$  can be expressed in terms of itself. An expression known as the *Bellman Equation* (ref).

$$Q(s, a) = \mathbb{E}[R_{t+1} + \gamma \sum_{a'} Q(s', a')] \quad (3)$$

The above form is handy because it opens doors for iterative approaches such as dynamic programming (ref).

## 2.4 Function Approximation

If one can calculate the Q-table from equation 3, they can successfully build an agent that maximizes the total reward. As mentioned in the last section, this can be easily done with dynamic programming. However, as the number of states grows largers, dynamic programming and other iterative approaches become unfeasible due to current memory and time limitations. Fortunately, it turns out that the Q-table can be approximated instead of having its exact values determined, and it still produces great results. This work tries to achieve that using linear regression and deep neural networks.

## 2.5 Temporal Difference Learning and Monte Carlo

In Monte Carlo Approaches, the agent plays an entire episode, keeping track of the rewards received at each timestep. After that, it updates the value function for each visited state based on the following equation.

$$V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)] \quad (4)$$

Where  $\alpha$  is the learning rate (ref).

Therefore, the agent only learns after an entire episode has been played.

In Temporal Difference Learning, however, the value of  $V$  is updated after each timestep. At time  $t + 1$ , the observations made during time  $t$  are already being considered. In its simplest form, the method is called TD(0) or one step TD, and its update equation is as follows.

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)] \quad (5)$$

The previous expression is referred as the *TD(0) error*.

## 2.6 Q-learning

Q-learning is another technique based on Temporal Difference Learning to learn the Q-table. The main difference between it and the previous shown technique TD(0) is that Q-learning is off-policy, and TD(0) is on-policy.(ref) This is reflected in its update equation, which is as follows.

$$Q(s, a) \leftarrow Q(s, a) + \alpha[R_{t+1} + \gamma \max_{a'} Q'(s', a') - Q(s, a)] \quad (6)$$

The fact that there is no constraint regarding the action  $a'$ , only that it must optimize  $Q'$ , makes it an off-policy method.

### 2.6.1 Deep Q-learning

In order to approximate the Q-table to make it feasible even when the number of states is large, since Google AlphaGo's paper (ref), it has becoming common the use of a deep neural network.

This work also makes use of that.

## 2.7 TD-Lambda

TD(0) is biased as it seems information from a single timestep in order to perform an update. It does not take into account the fact that the action that caused a reward might have happened several timesteps earlier, which can lead to slow convergence. Monte Carlo methods, although not biased, have a lot of variance since they use the rewards of an entire episode to perform its update.

From equation 4, we can define the 1-step return,  $G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$ . We can extend the concept to 2-step return,  $G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})$ , and, generically, to,  $G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n V(S_{t+n})$ .

TD-Lambda methods use a mathematical trick to average all the possible n-step returns into a single one. This is done by introducing a factor  $\lambda \in [0, 1]$  and weighting the nth-return with  $\gamma^{n-1}$ . It can be shown (ref) that when  $\lambda = 0$ , the method is equivalent to TD(0), and when  $\lambda = 1$ , equivalent to Monte Carlo. So, intuitively, by setting  $0 < \lambda < 1$ , we can get a mixture of both methods (ref).

### 3 Experiment

#### 3.1 Modeling

*TODO*

#### 3.2 Results and Discussion

*TODO*

### References

- [1] T. Guardian. Alphago seals 41 victory over go grandmaster lee sedol. <https://www.theguardian.com/technology/2016/mar/15/googles-alphago-seals-4-1-victory-over-grandmaster-lee-sedol>, 2016. accessed June 05, 2018.