



THE UNIVERSITY OF EDINBURGH

INFORMATICS HONORS PROJECT

Generative Adversarial Networks Generating Novel Reinforcement Learning Policies

Author:
Giovanni ALCANTARA

Supervisor:
Dr. Timothy HOSPEDALES

*A thesis submitted in fulfilment of the requirements
for the degree of Bachelor's of Engineering*

in the

School of Informatics
College of Science & Engineering

March 24, 2018

THE UNIVERSITY OF EDINBURGH

Abstract

Faculty Name
College of Science & Engineering

Bachelor's of Engineering

**Generative Adversarial Networks Generating Novel Reinforcement Learning
Policies**

by Giovanni ALCANTARA

TODO:

Contents

Abstract	iii
1 Introduction	1
1.1 Motivation	1
1.2 Goal	2
1.3 Data pipeline	2
1.4 Structure of the report	3
1.5 Main contributions	3
2 Background	5
2.1 Reinforcement Learning	5
2.1.1 Markov Decision Processes	5
2.1.2 Q-learning	5
2.1.3 Exploration/Exploitation tradeoff	6
2.1.4 Deep Reinforcement Learning	7
DQN	7
2.1.5 Problems with reinforcement learning techniques	7
2.1.6 Actor critic	7
2.2 Generative Adversarial Networks	7
2.2.1 Architecture of GANs	7
2.2.2 Successes	8
DCGAN	8
2.2.3 Conditional GANs	8
2.3 Existing related work	8
2.3.1 Generative Adversarial Imitation Learning	8
3 Environment	9
3.1 OpenAI Gym	10
3.1.1 Motivation	10
3.1.2 Algorithmic environments	11
3.1.3 MuJoCo and physics environments	11
3.1.4 Other environments	11
3.2 Baseline: FrozenLake-v0	11
3.2.1 Description of the task	11
3.2.2 Motivation and shortcomings	12
3.3 Extended baseline: Randomised Frozen Lake	12
3.3.1 Adjusting the reward system	15
4 Dataset creation	17
4.1 Q-learning on RandomisedFrozenLake	17
4.2 Experiment set up	17
4.3 Distributed Q-learning	19
4.4 Analysis of results	20

4.5	Transferable knowledge	20
5	Adversarial Networks Training	23
5.1	Policy GAN	23
5.1.1	Gradient descent optimisation	24
5.1.2	Hidden layers	25
5.2	Generator	25
5.2.1	Activation functions	25
5.2.2	Batch Normalisation	25
5.3	Discriminator	26
6	Generative Adversarial Q-learning	29
6.1	Using the trained Discriminator	29
6.1.1	Speeding up Q-learning on unseen maps	29
6.2	Using the trained Generator	29
6.2.1	Initialisation	29
6.2.2	Exploration	29
7	Scaling up to more complex tasks	31
7.1	Larger Frozen Lakes	31
7.2	Physics environments	31
8	Conclusion	33
	Bibliography	35

Chapter 1

Introduction

1.1 Motivation

Reinforcement learning is the problem faced by an agent that must learn behaviour through trial-and-error interactions with an environment.

Reinforcement learning as a field has had major successes in the past few years (Tesauro, 1995; Singh et al., 2002; Kohl and Stone, 2004; Ng et al., 2006), particularly as techniques utilising deep neural networks (DNN) have started permeating the research community. Techniques like Q-network (Mnih et al., 2015), trust region policy optimisation (Schulman et al., 2015), and asynchronous advantage actor-critic (A3C) (Mnih et al., 2016) helped stem an area of research of recent significant importance: deep reinforcement learning (DRL) (Arulkumaran et al., 2017).

Traditional reinforcement learning approaches generally lacked scalability, limiting these techniques to fairly low-dimensional problems. These limitations are in terms of memory complexity and computational complexity (Kaelbling, Littman, and Moore, 1996).

As such, using these usually becomes intractable when modelling real-world systems, due to the many variables and unknowns that are present in such systems (Strehl et al., 2006).

As DRL methods rectified some of these issues, new ones started to emerge, particularly limitations inherent to using deep neural networks. Notably, the need to have access to large datasets for training, particularly if in the context of applications that require image processing, such as autonomous vehicle control (Krizhevsky, Sutskever, and Hinton, 2012), has proven to be a critical limitation.

In real-world reinforcement learning applications, environment observations often rely heavily on computer vision and image processing (Berns, Dillmann, and Zachmann, 1992), which often provide an incomplete picture of the state that the agent is in. In such types of scenario, formally known as partially observable Markov decision processes (POMDP) (Monahan, 1982), not only do we have fragmentary observations, but it is also sometimes prohibitive to build large datasets that DRL requires to train the agent.

In deep learning, one of the ways to circumvent this constraint is *transfer learning*, the ability to leverage models trained in a particular domain on different applications. Transfer learning has proven pivotal in achieving successes in a wide variety of applications, without the need to train expensive models from scratch (Pan and Yang, 2010).

There has been much work in improving transferability of reinforcement learning models, most notably multitask learning (Caruana, 1998) and curriculum learning (Bengio et al., 2009), though few methods have tried to bootstrap learning with the use of DNNs.

There are few successes in transferring deep reinforcement learning across domains. Mostly notably, Jaderberg et al., 2016 introduced a technique to identify, in an unsupervised way, multiple pseudo-reward functions based on all training signals that the agent collected as observations. While doing deep reinforcement learning, therefore, Jaderberg et al. would not only try to directly maximise the agent’s cumulative reward, but also all the identified extrinsic rewards. There is a potential to use these identified extrinsic rewards in other domains, but this is a backwards way to tackle the problem of transferring behaviour. Here, we would first identify what auxiliary rewards the observations can give, to then reuse them on different task. Also, we have little to no control to guide the unsupervised exploration of auxiliary rewards functions towards a related task that we have the power to define. In fairness, Jaderberg et al. introduced auxiliary rewards as a way to speed up reinforcement learning on a single task, rather than aiming to transfer these to related tasks.

In our work we provide a general framework that lets us speed up reinforcement learning on unseen tasks in related domains. We do this by training deep learning models over a distribution of optimal policies for different configurations of a task in a certain domain. More specifically, given a distribution of trained policies in variations of an environment, we train two models: a generative model that is able to generate policies for different configurations of a task in a domain, and a discriminative model that is able to tell whether a policy is a good one within this domain. We show how using these models while doing reinforcement learning can speed up learning on new unseen configurations.

The generative model and the discriminative model are trained using deep neural networks in an adversarial architecture also known as a Generative Adversarial Network (GAN), introduced in Goodfellow et al., 2014’s seminal work. While this idea was popularised with applications in image synthesis, most notably Deep Convolutional Generative Adversarial Networks (Radford, Metz, and Chintala, 2015), there have been successes using GANs within reinforcement learning.

Specifically, work on generative adversarial imitation learning (Ho and Ermon, 2016) has shown remarkable speedups in the task of imitating behaviour given expert policies.

1.2 Goal

Given all these motivations, we can now define the **final goal of our work**:

Transferring knowledge obtained with reinforcement learning techniques on certain tasks to similar, unseen tasks, by leveraging the power of Generative Adversarial Networks (GANs). In this way, we hope to achieve shorter training times and better rewards when training models for these unseen tasks.

1.3 Data pipeline

In this section we provide a bird’s eye view to the whole project by showing a visual schematic of the data pipeline (Figure 1.1)

We will be referencing this pipeline schematic throughout the whole report. In fact, as we dive deeper into technical details of each component of the data pipeline,

it is critical to take them in the context of the whole pipeline as to justify certain design decisions we make for our experiments.

At this point, it is important to highlight how we achieve the goal we set in the previous subsection (subsection 1.2) with this data pipeline.

The input to our system is a set of similar tasks that would be traditionally solved with reinforcement learning techniques. What makes these tasks similar is the fact that they where the goal is the same, but the agent operates in different map configurations. We provide more details on the environment in Chapter 3. We split these tasks into a test set and a training set. We train all of these maps configurations with Q-learning, a reinforcement learning technique to obtain optimal policies.

Now, what we want to achieve, as per our goal, is:

1. train a GAN system to capture the distribution of the policies in the training set,
2. use the trained GAN to develop a new algorithm that is able to transfer that knowledge to speed up training of tasks in the test set of map configurations,
3. verify and benchmark that the test policies obtained with the new method are better than the ones trained with traditional Q-learning.

1.4 Structure of the report

Our report is structured in a what that strongly follows each component of the data pipeline outlined in the previous subsection.

- Chapter 2 introduces important background information on reinforcement learning, Markov Decision Processes, Q-learning and Generative Adversarial Networks.
- Chapter 3 details the environment being used in our experiments.
- Chapter 4 reports on the process of training policies on all the different tasks with traditional reinforcement learning (Q-learning).
- Chapter 5 outlines the process of training the GAN given the trained policies in the training set.
- Chapter 6 describes the new algorithm we develop to speed up learning on unseen tasks.
- Chapter

1.5 Main contributions

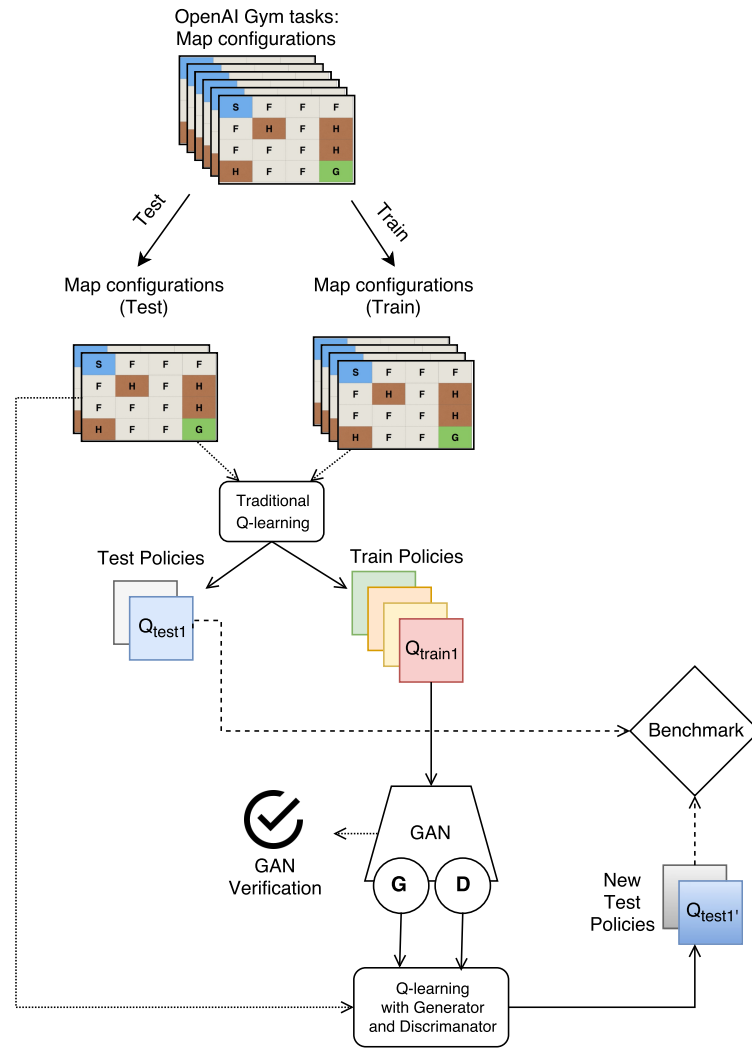


FIGURE 1.1: The Data Pipeline

Chapter 2

Background

2.1 Reinforcement Learning

2.1.1 Markov Decision Processes

Environments in traditional reinforcement learning application are usually modelled as Markov Decision Processes or MDPs.

These can be formulated as systems with the following components:

- Finite set of states $S = \{s_0, \dots, s_n\}$ and actions $A = \{a_0, \dots, a_m\}$.
- A distribution of probabilities $P_a(s, s')$ for transitions from state s to s' for each possible action a .
- A reward function $R : S \mapsto \mathbb{R}$ for being at a particular state.
- The goal in reinforcement learning is to maximise the final reward that an agent achieves in the environment.

As the name suggests, MDPs obey the *Markov property*, whereby the probability of the system being in a certain future state exclusively depends upon the present state, and not upon an arbitrarily-long sequence of past states.

In figure 2.1 we show a sample schematic of a Markov Decision Process, and how states, actions, and rewards could be connected between each other.

2.1.2 Q-learning

Now that we contextualised reinforcement learning environments as Markov Decision Processes, we can introduce the final objective of reinforcement learning tasks: finding a function $\Pi : S \mapsto A$ called **policy**, that maps the appropriate action $a \in A$ given the current state $s \in S$, as to maximise our agent's final reward.

In particular, we will now present Q-learning (Watkins and Dayan, 1992), an algorithm that yields an optimal policy given an MDP.

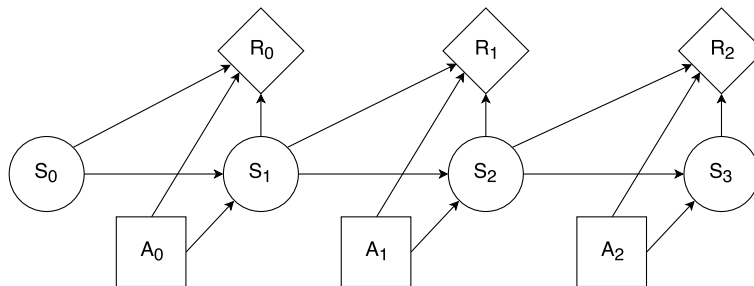


FIGURE 2.1: Sample schematic of an MDP.

Q-learning lets us learn the *quality*, or expected utility, for each state-action combination. That is, for each state, let's estimate all the expected rewards we obtain by taking each possible action at that particular state.

More formally, we estimate a function $Q : S \times A \rightarrow \mathbb{R}$. We can model Q as a mapping table (initialised with some uniform values), whose value we update at each time step of our simulations.

Here's how we update our Q-table at each time step t :

$$\begin{aligned}
 Q(s_t, a_t) &= \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \times \left[\underbrace{r_{t+1}}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right] \\
 Q_{t+1} &= Q_t - \underbrace{\beta}_{\text{discriminator learning rate}} \times \underbrace{\nabla D(Q_t)}_{\text{gradient of D at } Q_t} \\
 Q(s_t, a_t) &= \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \times \left[\underbrace{r_{t+1}}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right] \\
 &\quad - \underbrace{\beta}_{\text{discriminator learning rate}} \times \underbrace{\nabla D(Q)(s_t, a_t)}_{\text{gradient of D at } Q_t} \\
 &\quad \frac{\sum_i^{\text{test set}} \text{valid_score}(Q_i) - \text{valid_score}(Q'_i)}{\#(\text{test set})}
 \end{aligned}$$

where:

- $\alpha \in [0, 1]$ is the *learning rate*, a coefficient that regulates how much the newly learned values will contribute in the update
- $\gamma \in [0, 1]$ is the *discount factor*, a coefficient that controls the weight of future rewards. Values closer to 0 will make our agent "short-sighted", considering only the immediate rewards.

What is our optimal policy when we do Q-learning then? After training, it is simply that function $\pi : S \rightarrow A$ that, for each state, returns the action with maximum expected utility in our Q-table.

2.1.3 Exploration/Exploitation tradeoff

An important theme in reinforcement learning, and that we heavily focus our attention on in our project is the idea of the tradeoff between Exploration and Exploitation.

2.1.4 Deep Reinforcement Learning

DQN

2.1.5 Problems with reinforcement learning techniques

2.1.6 Actor critic

2.2 Generative Adversarial Networks

Generative Adversarial Networks, or GANs, are a deep neural network architecture composed of two neural networks, set against each other (thus "adversarial").

GANs were introduced in a paper by Goodfellow et al., 2014. Yann LeCun, Facebook's AI Research director referred GANs as being "the most interesting idea in the last 10 years in ML".

What have made GANs particularly successful in past work is their ability to model any distribution of data, and provide, after training, both a generative model and a discriminative model based on the initial distribution of the data. Let us define what we mean by these terms.

A *generative* model is one that is able to produce new data points that fit the distribution of the training data. For example, a Gaussian Mixture Model is a model that, after training, could generate new data based on the original distribution. More formally, training a generative model can be seen as maximum likelihood estimation problem:

$$\mathcal{L} = \prod_{i=1}^m p_{\text{model}}(x_i; \theta)$$

More specifically, we want to maximise the likelihood that a sample x that the generative model produces belongs to the distribution of our initial data.

A *discriminative* model is one that can discern (or "discriminate") the difference between two (or more) classes/labels. For example, we could train a Convolutional Neural Network that is able to tell us whether an input image is a face (1) or not (0).

Depending on the task at hand, we may be more interested in the generative model (for example in the case of image synthesis), or in the discriminative model, or in both.

Next, we will look at the typical architecture of Generative Adversarial Networks.

2.2.1 Architecture of GANs

Since GANs were introduced in 2014, many variants of GANs spawned from the research community. Here we will focus on the architecture of the original, vanilla GAN.

All variations of GANs will nevertheless contain two main components, which are both modelled as neural networks. These two models are trained simultaneously. One is the *generator* model, that takes an input z and produces a sample x from an implicit probability distribution. The other is the *discriminator*, a classifier that, given a sample, tries to identify it as originating from the generator model or from the original distribution.

An intuition to this set up is the following: the generator G can be imagined as a counterfeiter trying to produce fake banknotes or paper money; the discriminator D is the police, trying to distinguish real banknotes from the fake ones generated

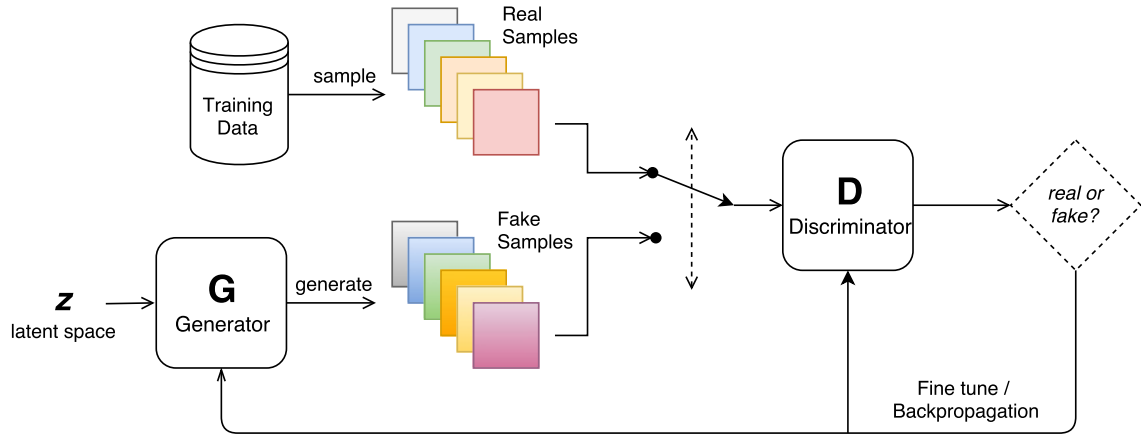


FIGURE 2.2: Architecture of a Generative Adversarial Network

by the counterfeiter. As the counterfeiter becomes better and better at producing banknotes, the police will also try to improve its counterfeiting techniques.

More formally, G is a differentiable function whose parameters are trained to minimise correct assignments of D . D is also a differentiable function, which has been trained to maximise correct labels to the real and the fake samples.

In figure 2.2 summarises this whole procedure by showing the architecture of GANs.

At each iteration, the mini-batch inputs to the discriminator D are taken from both the real data sampled from the training data, and the samples generated by the generator, whose inputs is random latent variable z . Given D 's prediction, we then apply a gradient-based optimisation method to update both D 's and G 's parameter.

To do so, we follow the loss function that defined based on the goals of each network. This cost function implicitly models a zero-sum game, which in game theory are guaranteed to achieve an equilibrium, by the minimax theorem (Du and Pardalos, 2013).

The cost function for the whole system would look like the following minimax game:

$$\min_G \max_D V(G, D) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_g(z)} [\log(1 - D(G(z)))]$$

2.2.2 Successes

DCGAN

2.2.3 Conditional GANs

2.3 Existing related work

2.3.1 Generative Adversarial Imitation Learning

Chapter 3

Environment

In this chapter we report the process that went into choosing the environments that we will be using in the project's simulations, and from which we will be basing our simulations.

There are many choices of environments and tasks that are publicly available. Some of these we will be introducing in this chapter.

What makes this step non-trivial and deserving of its own chapter is that different reinforcement learning techniques are more suitable to different categories of tasks. Similarly, different machine learning and deep learning techniques are more or less efficient when applied to different tasks.

In a research effort that is heavily dependent on building reinforcement learning and deep learning models, the choice of environment is a critical one.

Furthermore, we also need to achieve this without losing focus on the main motivations for the whole project (Chapter 1): *optimising reinforcement learning algorithms by adding transferability of pre-trained models on unseen maps or configurations of a task*.

This last point implies that a substantial part of the computational work in the project will be about training hundreds of thousands of reinforcement learning models to build a dataset over a distribution of different maps (we present this in Chapter 4). This is an important point: in the many experiments we ran, this turned out to be the biggest bottleneck, and required plans to distribute computations across multiple machines to make the computational time feasible in the timespan allocated to the project.

With these points in mind and given the experimental nature of the work, we conclude that a bottom-up approach in complexity is preferable. Given successful results with "easier" tasks, we can scale up in complexity and hopefully formalise and generalise our approach to more tasks (Chapter 7).

Easier tasks will enable us to explore different reinforcement learning approaches that we introduced in Subsection 2.1 with the guarantee that they will give satisfiable results. We can use these results as a foundation of the further steps (specifically Generative Adversarial Networks training in Chapter 5).

Before introducing candidate environments, let us define what is meant by an "easy" reinforcement learning task. What we are looking for is ideally a task with a discrete and relatively small set of observable states and actions. Why does this condition make the task easier?

Imagine building a decision tree (such as the one shown in figure 3.1 with all possible states-action transitions, until we either: 1) reach a goal state, or 2) reach an arbitrarily maximum iteration time step $t = \eta$ or depth of the tree (to prevent infinite iterations). Also assume we were building this tree in a brute-force manner (worst-case scenario of reinforcement learning resolutions), such that we need to build all possible paths or trajectory that the agent will need to take. Therefore, we would

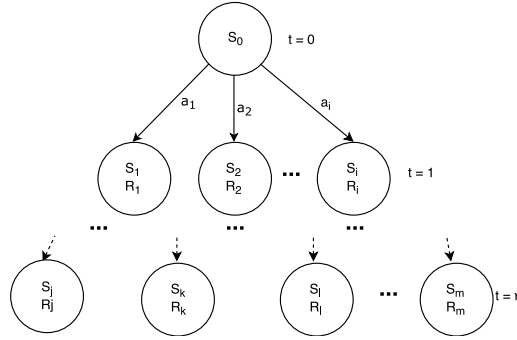


FIGURE 3.1: Sample state-action tree

need to visit each node of the tree, until we arrive at the leaves. At the leaves, we can find the achieved reward for the agent given the path it has taken to get there.

The breath and depth of the state-action tree will increase as we increase the possible set of states we would need to traverse, adding up to the space and time complexity of our solution, which is exponential in $\#S$ and $\#A$ for both space and time ($\#S$ indicates the cardinality of a set S).

With that said, there is a solid line between our ideal environment and a trivial one that could be solved without the aid of reinforcement learning techniques.

Let's explore some candidate environments next.

3.1 OpenAI Gym

OpenAI Gym (Brockman et al., 2016) provides a toolkit that aids in the process of building reinforcement learning systems and evaluating algorithms to solve such tasks.

OpenAI Gym provides an environment interface `Env`. The interface abstracts the following operations on an environment:

- `step(action)` – Simulates one time step by executing the provided action. This returns observation (the new current state), the reward for taking such action, a flag `done` indicating whether the system has reached a final state, and `info` providing additional information dependent on the environment.
- `reset()` – Resets the environment, i.e. the initial state is restored
- `render()` – Renders the environment in human-readable format.

Now, we could either build implementations for such interface (if we were to implement our own environment's logic), or use the provided implementations for several environments in the OpenAI Gym library, which includes board games, algorithm-based or physics-based environments, Atari games, etc.

3.1.1 Motivation

In this project we will mostly deal with environments provided in the OpenAI Gym. There are several reasons for such decision:

- It abstracts the need to implement the logic of a separate environment. Implementing our own environment adds a point of failure to our whole (quite

experimental) work, as well as imposing a bigger time constraint to the one we already have;

- OpenAI Gym has become a standard academic tool for reinforcement learning researchers, therefore many papers and articles build on top of this framework;
- Environment implementations in the OpenAI Gym library are constantly expanding and being revised by an active community, also thanks to the support of the overarching organisation (OpenAI);
- The core implementation of the gym module ([openai/gym on Github](#)) allows for straightforward extensions on existing environments, which, as we will see in subsection 3.3, will be critical in this project;
- While sadly not applicable anymore, OpenAI Gym used to provide and support an online platform for developers to compare the performance of different reinforcement learning algorithms for each task. This was in a form of a leaderboard, measuring the performances, as well as providing the implementation and data of winning techniques. Such data could have been used directly as input to our Generative Adversarial Networks.

3.1.2 Algorithmic environments

3.1.3 MuJoCo and physics environments

3.1.4 Other environments

3.2 Baseline: FrozenLake-v0

One of our baseline environments is FrozenLake-v0 ([OpenAI Gym](#)), one of the algorithmic environments provided in the OpenAI Gym. In this section we describe the task, motivation for choosing it as one of our baseline models, and some shortcomings that we faced in using this environment in our project's pipeline.

3.2.1 Description of the task

In FrozenLake-v0 we control an agent in a grid world, more precisely the grid shown in figure 3.1. The objective is to go from a starting tile (S) to a goal state (G), by moving in four possible directions from a tile: up, down, left and right.

What differentiates this from a trivial path-finding or planning problem? A couple of things:

1. there are both walkable and non-walkable tiles in the grid (these are respectively frozen tiles F and holes H), and
2. tiles are "slippery", as in the agent could "slip" while navigating the grid world, meaning that the movement direction of the agent is uncertain and only partially depends on the direction we tell the agent to follow, i.e. the direction is non-deterministic.

The reward in FrozenLake-v0 is 1 for reaching the goal state, and 0 for being in any other state. The system stops executing when the agent falls into a hole, and the environment needs to be reset.

S	F	F	F
F	H	F	H
F	F	F	H
H	F	F	G

TABLE 3.1: FrozenLake-v0's default 4x4 configuration

3.2.2 Motivation and shortcomings

FrozenLake-v0 is generally classified as a straightforward reinforcement learning task. An optimal solution could be found by creating a model of the environment, that is merely recording where frozen tiles are while exploring the grid world. This is a *model-based* approach to reinforcement learning, and it is perhaps less interesting than learning by exploration, without being biased by the particular configurations of the map.

Model-free algorithms, like Q-learning, need no accurate representation specific to the environment, and they are therefore more transferable to different configurations of the map or even to different tasks, which is the final objective of our project.

In fact, let's take the case of Q-learning applied to FrozenLake-v0: we do not need to have an explicit knowledge of the dynamics of each different tiles. We do not build a policy that explicitly favours movements towards frozen tiles or towards the goal. In fact, the agent does not have a model of what a frozen tile is, nor a model of the goal tile or a hole—it just learns by exploration that it is rewarded when it gets to the goal state, and that is what it implicitly aims for.

This is the sort of behaviour that we can transfer to different tasks—again our ultimate objective.

In its current OpenAI Gym implementation, FrozenLake-v0 is a static map with the fixed configuration that is shown in figure 3.1. There is currently no way to generate random configurations of the map, so next up, we will be extending this implementation to account for that.

3.3 Extended baseline: Randomised Frozen Lake

We would like to extend OpenAI Gym's implementation of FrozenLake-v0 so that it can generate random configurations of the map.

Before we move on, let us contextualise this in the bigger picture as to not lose focus of what we are trying to achieve. Why do we need different configurations of the map, again? We want to train reinforcement learning models on different configurations so that we can have a distribution of policies over different maps which we can use as input to our GAN. After training our GAN, we will have a Generator network spawning new policies for unseen configurations without having to find it through (computationally expensive) reinforcement learning algorithms!

Listing 3.1 shows how we implemented the algorithm to generate random maps. The critical line is line 4, which uses numpy's `random.choice()` method that samples a matrix of a given size, given probabilities for each element ('F' and 'H' tiles in our case).

So, we can pass in the desired size of the map. By default, it will generate 4x4 grids like the one in FrozenLake-v0, but we could generate maps of arbitrary size,

which will result in higher task complexity. We explore these harder extensions in Chapter 7.

We can also pass it the probability that a tile will be a frozen one through the parameter p . The presence of fewer frozen tiles, and therefore of more holes, makes the goal harder to achieve for the agent.

```

1 def generate(size=4, p=0.8):
2     valid_map = False
3     while not valid_map:
4         config = np.random.choice(['F', 'H'], (size, size), p=[p, 1-p])
5         config[0][0] = 'S' # set top left to be the starting tile
6         config[-1][-1] = 'G' # set bottom right to be goal tile
7         valid_map = is_valid(config)
8         p *= 1.05 # increase probability of frozen tile
9     return "".join(x) for x in config

```

LISTING 3.1: Algorithm to generate random configurations. Utility function `is_valid()` of line 7 is shown in listing 3.2

Notice how the `generate()` function in listing 3.1 only returns valid maps, that is, maps that have at least one frozen path from start to goal. Surely, we could train models on environment configurations that are not solvable. Q-learning, for example, would just return a Q matrix with all Q-values equals to 0, since the agent will never get to the goal tile and get its reward. If it is unclear why, refer to the Q-learning subsection 2.1.2.

Using such constraint on map validity, we can limit the number of models we need to train by a significant amount, therefore reducing training time. For the 4x4 grid, rather than having $2^{(16-2)} = 2^{14} = 16,384$ possible configurations, we can reduce it to only 3,827 possible valid map configurations.

To check whether a map is solvable we use depth-first search from the start tile to the goal. If there is such path, then it is a valid map. Listing 3.2 shows the algorithm:

```

1 def is_valid(arr, r=0, c=0):
2     if arr[r][c] == 'G':
3         return True
4
5     tmp = arr[r][c]
6     arr[r][c] = "#" # temporary mark with '#' to remember visited tiles
7
8     if r+1 < size and arr[r+1][c] not in '#H': # go down
9         if is_valid(arr, r+1, c) == True:
10             arr[r][c] = tmp
11             return True
12     if c+1 < size and arr[r][c+1] not in '#H': # go right
13         if is_valid(arr, r, c+1) == True:
14             arr[r][c] = tmp
15             return True
16     if r-1 >= 0 and arr[r-1][c] not in '#H': # go up
17         if is_valid(arr, r-1, c) == True:
18             arr[r][c] = tmp
19             return True
20     if c-1 >= 0 and arr[r][c-1] not in '#H': # go left
21         if is_valid(arr, r, c-1) == True:
22             arr[r][c] = tmp
23             return True
24
25     arr[r][c] = tmp
26     return False

```

LISTING 3.2: Depth-first search to check if a frozen lake map is valid

So the `generate()` function returns a valid random frozen lake map represented as a list of strings. An example of an output is `["SFHH", "HFHH", "HFHH", "HFFG"]`. Each string in the list encodes a row configuration of the frozen lake.

What is left to do now, is to extend `FrozenLake-v0` so that we could pass in any map configuration, therefore exposing our randomly generated maps to the gym environment interface described earlier in section 3.1. Listing 3.3 shows how to register a new environment by extending a pre-existing one.

```

1 random_map = generate(size=4, p=0.8)
2 register(
3     id='RandomisedFrozenLake',
4     entry_point='gym.envs.toy_text:FrozenLakeEnv',
5     kwargs={'is_slippery': True, 'desc': random_map},
6     max_episode_steps=100,
7 )

```

LISTING 3.3: Code to extend `FrozenLake-v0` with random maps.

Now, to start doing simulations on our new environment we can initialise it just like any other OpenAI Gym environment:

```
1 env = gym.make('RandomisedFrozenLake')
```

We created a pull request to the OpenAI Gym's Github repository integrating this feature supporting random maps for `FrozenLake-v0`, so that developers and researchers could make use of these functionalities (*Add ability to generate random frozen lake maps by gosi · Pull Request #835*).

3.3.1 Adjusting the reward system

The vanilla OpenAI Gym implementation of FrozenLake-v0 has sparse rewards, meaning that the agent is rewarded a value of 1 if it reaches the goal, and a value of 0 if it falls into a hole.

The problem with this system of rewards is the fact that it does not take into account the number of timesteps that it takes the agent to solve the task.

Consider the three paths shown in figure 3.2. All three paths would yield a final reward of 1 in the OpenAI Gym's implementation. But intuitively, in a real-world scenario, we'd like to have a path-finding robot that is able to navigate the environment with as few actions as possible.

More generally, reinforcement learning environments have rewards that are inversely proportional to the number of steps that it takes the agent to solve (Bellemare et al., 2013; Brockman et al., 2016; Zamora et al., 2016)

We wish to leverage this alternative reward system to train agents that can therefore reach the goal state as quickly as possible.

In terms of extending our RandomiseFrozenLake implementation, we can just update the final reward as

$$r = \begin{cases} 0 & \text{if agent fell into a hole} \\ 6/t & \text{if agent reaches goal state} \end{cases}$$

where t is the number of timesteps needed to reach the goal state.

If the agent reaches the goal state in the minimum number of timesteps (6), the reward would be 1, otherwise we add a penalty for any additional timestep required.

The rewards of each path in figure 3.2 will therefore be, from left to right, $r = 1$, $r = 1$ and $r = 6/14 = 0.43$.

In section 4.4 and 4.5, we will show how this reward system (favouring paths that are shorter) will help us define the type transferable knowledge that we wish to transfer across tasks.

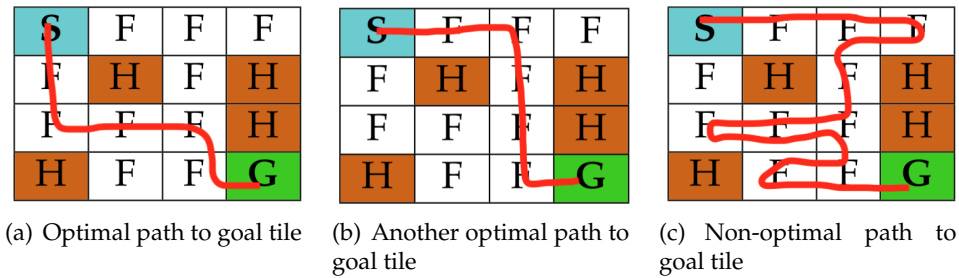


FIGURE 3.2: Different paths to the goal tile

Chapter 4

Dataset creation

In Chapter 3 we finalised our OpenAI environment choice and have a systematic way to generate different configurations of our environments.

The next step is to build up our dataset by solving as many of these configurations as possible with the use of reinforcement learning algorithms.

We justified in previous sections (specifically section 2.1.2 and section 3.2.2) how using *model-free* reinforcement learning algorithms will put us in the right direction to achieve the project's goals highlighted in section 1.1: improve the transferability of pre-trained models to different configurations, environments, and tasks.

One choice that is left to make before we move on to train models using our randomised environments is what reinforcement learning algorithm we should use.

4.1 Q-learning on RandomisedFrozenLake

We have presented Q-learning in section 2.1.2 as a method that generates a policy by building up a table Q of values corresponding to the expected utility of taking each action at each observable state.

This table Q is represented as a 2D-matrix of size $(\text{env.observation_space.n}, \text{env.action_space.n})$, which in the case of RandomisedFrozenLake is just a (16×4) matrix.

This is good news. Having a policy that is encoded as a 2-dimensional data structure makes it an obvious input to our Generative Adversarial Network later on in the next chapter. In fact, as we presented in section 2.2.2, GANs have proven successful in image synthesis applications, where inputs were images, that is 2D matrices.

The choice of Q-learning as our reinforcement learning technique therefore becomes preferable. Other model-free algorithms like policy search may not have a clear 2D representation in the way the trained parameter set θ is encoded.

4.2 Experiment set up

Before we proceed, we need to be able to record some details about the process of training our data. If our ultimate objective is to benchmark performance of traditional reinforcement learning against our proposed approach, we need to record running time of our training, and see if that improves at the end.

Let's encapsulate training of a configuration in an `Experiment` class, whose code is shown in listing 4.1. An `Experiment` is initialised with an OpenAI Gym environment, which in our case is an instance of a `RandomisedFrozenLake`, and an integer value `num_episodes`, indicating the number of independent simulations the Q-learning algorithm will be doing. By calling `Experiment`'s `run()` instance method, we actually start training the model, with `Q` being updated at each iteration.

The instance variable `score` could be used as an evaluation criteria of the effectiveness of our training. It records the average reward the agent achieves for each episode during training.

`Experiment` also has a utility method called `dumps()` that serialises all this data and allows us to save it on disk.

```

1 class Experiment(object):
2     def __init__(self, env, num_episodes=10000):
3         self.env = env
4         self.Q = np.zeros([self.env.observation_space.n, self.env.
5             action_space.n])
6         self.num_episodes = num_episodes
7         self.score = None
8         self.valid_score = None
9         self.start = None
10        self.end = None
11
12    def run(self):
13        self.start = datetime.now()
14        # -----
15        # Run Q-learning algorithm, saving the rewards of each episode
16        # ...
17        # ...
18        self.end = datetime.now()
19        self.score = sum(rewards)/self.num_episodes
20
21    def dumps(self):
22        return dumps({'Q': self.Q, 'start': self.start, 'end': self.end, '
23            score': self.score, 'num_episodes': self.num_episodes})

```

LISTING 4.1: Experiment wrapper class to train one instance of `RandomisedFrozenLake`

It's critical to be able to evaluate the quality of the Q-table after training. To do so we just use the optimal policy (that is the policy that picks the action that has the maximum expected utility according to the Q-table) and run it for a certain number of episodes. A validation score could be then defined by the average reward achieved at each episode. Listing 4.2 shows the code to achieve that.

```

1 def validate(self):
2     rewards = 0
3     for i in tqdm(range(num_episodes)):
4         while j < 200: # Limit to 200 time steps
5             j+=1
6             # Choose an action by pick best action from Q-table
7             a = np.argmax(Q[s,:])
8
9             # Get new state and reward from environment
10            s1,r,d,_ = env.step(a)
11            rewards += r
12            s = s1
13            if d == True:
14                break
15            self.valid_score = rewards / num_episodes

```

LISTING 4.2: Code to validate a trained Q-table

4.3 Distributed Q-learning

Now that we formalised our experiment setup, we can run Q-learning on each of the map configurations of our RandomisedFrozenLake. For the 4x4 grid there are 3,827 possible valid map configurations. It takes an average of 15 seconds to train each Q-learning table on a 2.5 GHz Intel Core i7 processor. On a single machine, it would take around 16 CPU hours to run the whole set of experiments.

To make this step of our pipeline faster and scalable to more data, we decide to set up a distributed processing on a cluster with multiple machines. More precisely, we set up a MapReduce framework (Dean and Ghemawat, 2004) implementation running on an Hadoop cluster (Shvachko et al., 2010) of 25 machines.

A MapReduce program is composed of a Map procedure that takes in some (large) input and performs a particular operation whose output is then fed into another Reduce procedure, which outputs the final result. The power of MapReduce is that the framework orchestrates the processing of these procedures by marshalling the distributed servers, running the various tasks in parallel, managing all communications and data transfers between the various parts of the system, and providing for redundancy and fault tolerance.

Figure 4.1 shows the distributed architecture to train multiple Q-learning instances. The input to the system is a list of strings, each representing the map configuration of a RandomisedFrozenLake, e.g. "SFHHHFHHHFHHHFFG". These inputs are fed into Mapper programs running on different machines. The mapper's task is to initialise the experiment and run Q-learning on the environment. It outputs a key-value pair, where the key is the string that uniquely identifies a map configuration, and the output is the Experiment object that encapsulates the already trained Q-table. We have yet to assign a validation score to this table we trained, and that's the job of the reducer.

Each experiment's result is then written to an output file by the reducer. Each line will be again a key-value pair, with the map configuration string as key and the Experiment result dumped in string format. In the following sections, we can just parse and load these results to conduct our further analysis.

While this particular architecture does not make full use of the power of MapReduce (combining, partitioning and sorting), it is an optimal and convenient pipeline to distribute our computations across a cluster of machines, thereby drastically reducing the training time of the whole experiment set.

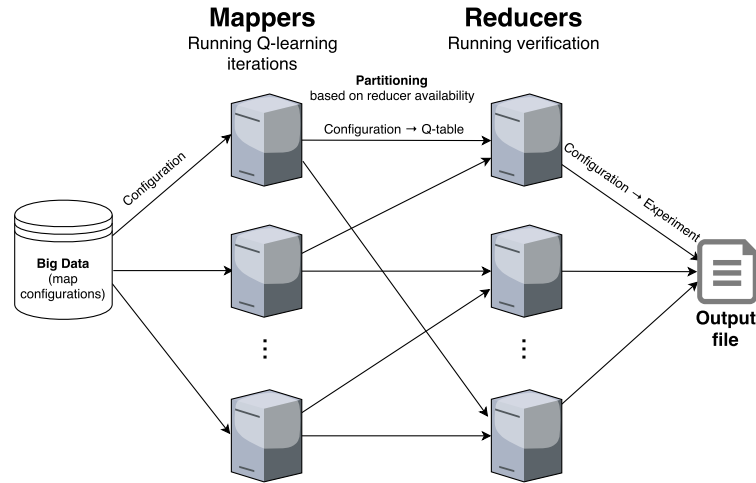


FIGURE 4.1: Schematic of distributed Q-learning with validation on MapReduce

4.4 Analysis of results

In this section we report some statistics on the data we just trained. This analysis will be useful when we perform our final benchmarking.

Figure 4.2 shows the average intensity values of the Q-tables in the training test.

We notice that the utility values in the second and third columns are, in most cases, the highest values at each row. We interpret this as our models favouring, on average, the action of going south (second column) and east (third column) over going west or north.

Indeed, according to the reward system that we described in section 3.3.1, our trained models will try and reach the goal state with as few actions as possible. This validates why on average, the best actions to make at each state will be going south or going east.

4.5 Transferable knowledge

Before proceeding, we need to start a discussion on the type of knowledge that we wish or we expect to be transferred. Having transferred knowledge is what will allow us to speed up training time on unseen tasks.

By using Q-learning to train policies on RandomisedFrozenLake, we effectively train agents to go from the top-left corner to the bottom-right corner. In section 3.3.1 we adjusted the rewards to

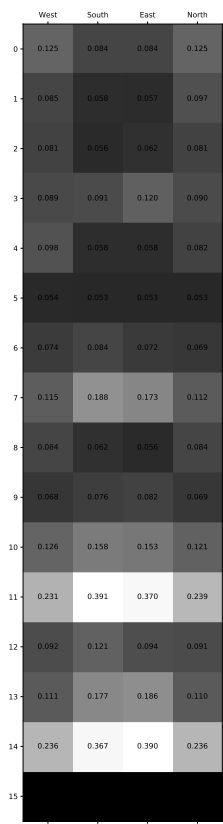


FIGURE 4.2: Schematic of distributed Q-learning with validation on MapReduce

Chapter 5

Adversarial Networks Training

Now that we built our dataset, we have the data we need to proceed with the adversarial training of our pipeline.

It is once again important to place this chapter in the context of our data pipeline, as we introduced it in section 1.3 and as we showed it in figure 1.1. The dataset we just created is composed of policies (which are just tables or 2D matrices). More specifically we have 3,827 such policies, corresponding to all the 3,827 different valid map configurations of our environment, `RandomisedFrozenLake`. We split the dataset in a training subset (80%) of the total number of policies, and the test subset (the remaining 20%).

In section 2.2 we introduced Generative Adversarial Networks, detailing their vanilla architecture, as they were introduced by Goodfellow et al. We now build off of this architecture to build a GAN architecture for policy generation.

In this section we present such architecture, explain the design decisions that were involved in configuring it, and describe the process of training it in an adversarial manner. We also provide a detailed breakdown of the generator and the discriminator's deep neural networks.

The input to our Policy GAN will only be the training subset of the policies we trained with Q-learning. In this way, we hope to be able to model, through our GAN, a distribution on these policies such that we can capture knowledge that we can transfer to other unseen tasks. These unseen tasks are exactly the ones corresponding to the map configurations of the test subset, and we therefore do not include those test policies as our input to our GAN.

5.1 Policy GAN

In figure 5.1 we show the architecture of our Policy GAN. We already described most of the details involved when training GANs in section 2.2.

Training a GAN is an iterative process that runs for a set amount of epochs (an epoch is one full training cycle on the training set). At each iteration, the mini-batch inputs to the discriminator D are taken from both the real data sampled from the training data (in our case policies from the training set that we created in chapter 4), and the samples generated by the generator. Given D 's prediction, we then apply a gradient-based optimisation method to update both D 's and G 's parameters.

Building neural net architecture that work well in practice generally involves finding and tuning many parameters and making several design decisions that make the issue not trivial. In the following subsections, we introduce such decisions and hyperparameters choice with the aim of allowing reproducibility of our results.

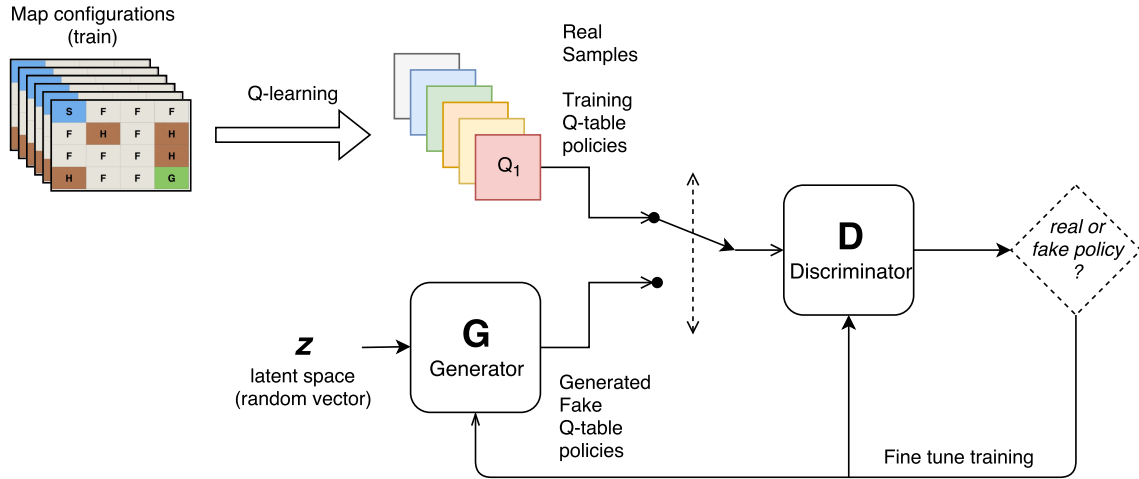


FIGURE 5.1: Our Policy GAN architecture

5.1.1 Gradient descent optimisation

Learning rates and gradient descent optimisations algorithms have historically been among the trickiest hyperparameters to set, as they drastically influence the final results, perhaps more than all other hyperparameters.

Gradient descent (Lemaréchal, 2012), a seminal approach to optimisation that dates back to Cauchy a few centuries back, still proves relatively successful in many machine learning applications. Since then, a lot of work has been put in devising algorithms that have adaptive learning rates and that lead to faster convergence. In our design, we specifically explored two such adaptive learning rates algorithms: RMSProp (Hinton, Srivastava, and Swersky, 2012) and Adam (Kingma and Ba, 2014)

Both perform local optimisation with different techniques and metrics that are constructed from the history of previous interactions.

RMSProp (algorithm 1) is a modification of the AdaGrad optimiser (Duchi, Hazan, and Singer, 2011), that modifies the gradient accumulation into a moving average that is exponentially weighted.

Adam (algorithm 2) is best described as a combination of RMSProp and Stochastic Gradient Descent (SGD) with momentum (Sutskever et al., 2013). Momentum accelerates SGD by multiplying the learning rate by a parameter that increases as we go towards the right direction in the gradient update.

Like in most applications involving deep learning, there is no definite way to see which approach would yield the best results without going through the whole training process.

We therefore empirically experiment with each of these approaches to see which choice is more suitable to our data distribution and yields the best results. "Best", in our case, is seen in the context of both the Generator network and the Discriminator network.

We look for an optimisation technique that leads to stable costs for both networks. In most deep learning problems, convergence of the cost function is limited to one single neural network. In adversarial architectures like GANs, or architectures with multiple deep neural networks like Variational Autoencoders (VAE), or Actor-Critic networks, convergences of the model is not a property that we can trivially evaluate (Kingma and Welling, 2013; Grondman et al., 2012).

Algorithm 1 RMSProp algorithm

Require: Global learning rate ϵ , decay rate ρ
Require: Initial parameter θ
Require: Small constant δ , usually 10^{-6} , used to stabilize division by small number
Initialize accumulation variables $r = 0$
while stopping criterion not met **do**
 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets y^i .
 Compute gradient: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$.
 Accumulate squared gradient:
 $r \leftarrow \rho r + (1 - \rho) g \odot g$
 Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + r}} \odot g$. ($\frac{1}{\sqrt{\delta + r}}$ applied element-wise)
 Apply update: $\theta \leftarrow \theta + \Delta \theta$
end while

We look for an asymptotic behaviour that has both networks converge to an optimum, without having one network prevail over the other. For example, we could have a Discriminator that produces perfect predictions after a few epochs of training, but that is a meaningless result if the Generator is not able to produce results that can "fool" the discriminative model.

The overall cost of the architecture, in short, must take into account both components:

$$\min_G \max_D V(G, D) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_g(z)} [\log(1 - D(G(z)))]$$

With both RMSProp and Adam we obtain much better convergence compared to vanilla SGD. While do not notice a noticeable difference between these two results, having adaptive learning rates clearly helps convergence for both models.

After optimising our hyperparameters with grid-search, our final optimal optimisation is the following:

Adam, step size $\epsilon = 0.0002$, exponential decay rates $\rho_1 = 0.5$, $\rho_2 = 0.999$.

5.1.2 Hidden layers

5.2 Generator

Figure 5.2 shows the

5.2.1 Activation functions

5.2.2 Batch Normalisation

Batch Normalisation (Ioffe and Szegedy, 2015) is a technique that improves stability and performance of neural networks by normalising the inputs of each layer such that they have a mean output activation of zero and standard deviation of one. Benefits of using batch normalisation include faster training time (i.e. faster convergence to optimal model). This also allows higher learning rates and weight initialisation won't be as critical.

Algorithm 2 Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)
Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$. (Suggested defaults: 0.9 and 0.999 respectively)
Require: Small constant δ used for numerical stabilisation (Suggested default: 10^{-8})
Require: Initial parameters θ
 Initialise 1st and 2nd moment variables $\mathbf{s} = \mathbf{0}, \mathbf{r} = \mathbf{0}$
 Initialise time step $t = 0$
while stopping criterion not met **do**
 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets y^i .
 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$.
 $t \leftarrow t + 1$
 Update biased first moment estimate: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$
 Update biased second moment estimate: $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$
 Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$
 Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$
 Compute update: $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$
 Apply update: $\theta \leftarrow \theta + \Delta \theta$
end while

We first compute the mean and variance of each hidden unit activation across the minibatch (size M):

$$\mu_i \leftarrow \frac{1}{M} \sum_{m=1}^M u_i^m$$

$$\sigma_i^2 \leftarrow \frac{1}{M} \sum_{m=1}^M (u_i^m - \mu_i)^2$$

The result of doing batch normalisation will then be:

$$u_i = w_i x$$

$$\hat{u}_i = \frac{u_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

$$z_i = \gamma_i \hat{u}_i + \beta_i = \text{batchNorm}(u_i)$$

where γ and β are parameters updated with gradient descent that will scale and shift the normalised activations.

5.3 Discriminator

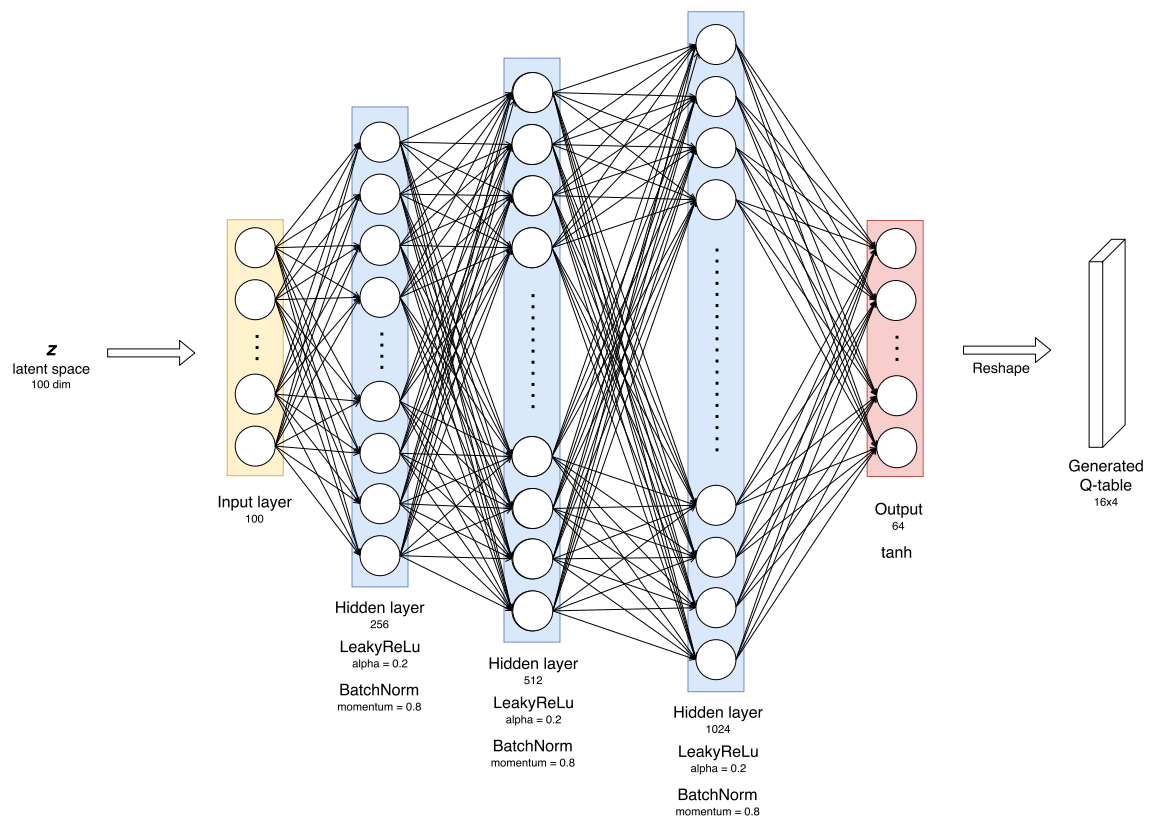


FIGURE 5.2: Architecture of the Generator network

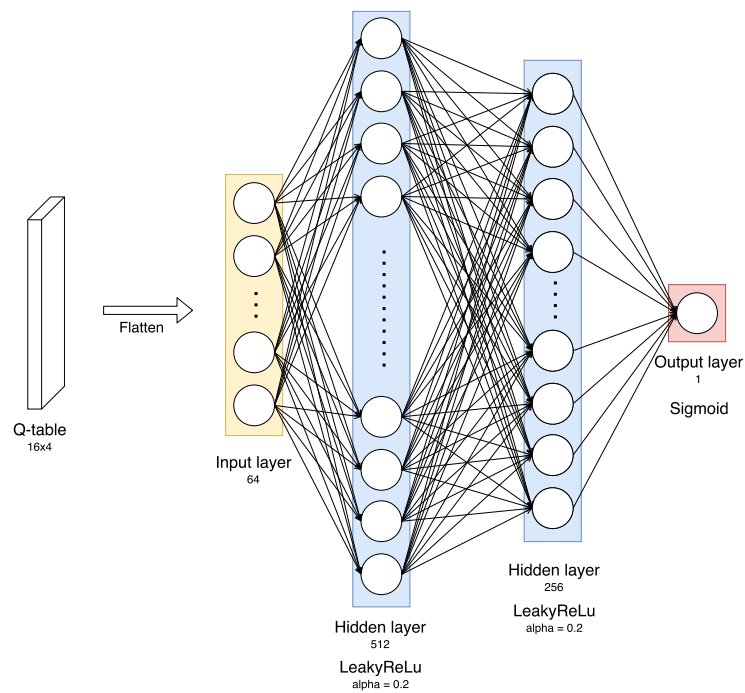


FIGURE 5.3: Architecture of the Discriminator network

Chapter 6

Generative Adversarial Q-learning

6.1 Using the trained Discriminator

6.1.1 Speeding up Q-learning on unseen maps

6.2 Using the trained Generator

6.2.1 Initialisation

6.2.2 Exploration

Chapter 7

Scaling up to more complex tasks

7.1 Larger Frozen Lakes

7.2 Physics environments

Chapter 8

Conclusion

Bibliography

- Add ability to generate random frozen lake maps by gvti · Pull Request #835. <https://github.com/openai/gym/pull/835>. (Accessed on 01/26/2018).
- Arulkumaran, Kai et al. (2017). "A Brief Survey of Deep Reinforcement Learning". In: CoRR abs/1708.05866. arXiv: 1708.05866. URL: <http://arxiv.org/abs/1708.05866>.
- Bellemare, Marc G et al. (2013). "The Arcade Learning Environment: An evaluation platform for general agents." In: *J. Artif. Intell. Res.(JAIR)* 47, pp. 253–279.
- Bengio, Yoshua et al. (2009). "Curriculum learning". In: *Proceedings of the 26th annual international conference on machine learning*. ACM, pp. 41–48.
- Berns, Karsten, Rüdiger Dillmann, and U Zachmann (1992). "Reinforcement-learning for the control of an autonomous mobile robot". In: *Intelligent Robots and Systems, 1992., Proceedings of the 1992 IEEE/RSJ International Conference on*. Vol. 3. IEEE, pp. 1808–1815.
- Brockman, Greg et al. (2016). *OpenAI Gym*. eprint: [arXiv:1606.01540](https://arxiv.org/abs/1606.01540).
- Caruana, Rich (1998). "Multitask learning". In: *Learning to learn*. Springer, pp. 95–133.
- Dean, Jeffrey and Sanjay Ghemawat (2004). "MapReduce: Simplified Data Processing on Large Clusters". In: *Proceedings of the 6th Conference on Symposium on Operating Systems Design & Implementation - Volume 6. OSDI'04*. San Francisco, CA: USENIX Association, pp. 10–10. URL: <http://dl.acm.org/citation.cfm?id=1251254.1251264>.
- Du, Ding-Zhu and Panos M Pardalos (2013). *Minimax and applications*. Vol. 4. Springer Science & Business Media.
- Duchi, John, Elad Hazan, and Yoram Singer (2011). "Adaptive subgradient methods for online learning and stochastic optimization". In: *Journal of Machine Learning Research* 12, Jul, pp. 2121–2159.
- Goodfellow, Ian et al. (2014). "Generative adversarial nets". In: *Advances in neural information processing systems*, pp. 2672–2680.
- Grondman, Ivo et al. (2012). "A survey of actor-critic reinforcement learning: Standard and natural policy gradients". In: *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 42.6, pp. 1291–1307.
- Hinton, Geoffrey, Nitish Srivastava, and Kevin Swersky (2012). *Neural networks for machine learning lecture 6a overview of mini-batch gradient descent*.
- Ho, Jonathan and Stefano Ermon (2016). "Generative adversarial imitation learning". In: *Advances in Neural Information Processing Systems*, pp. 4565–4573.
- Ioffe, Sergey and Christian Szegedy (2015). "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift". In: CoRR abs/1502.03167. arXiv: 1502.03167. URL: <http://arxiv.org/abs/1502.03167>.
- Jaderberg, Max et al. (2016). "Reinforcement learning with unsupervised auxiliary tasks". In: *arXiv preprint arXiv:1611.05397*.
- Kaelbling, Leslie Pack, Michael L Littman, and Andrew W Moore (1996). "Reinforcement learning: A survey". In: *Journal of artificial intelligence research* 4, pp. 237–285.

- Kingma, Diederik P. and Jimmy Ba (2014). "Adam: A Method for Stochastic Optimization". In: *CoRR abs/1412.6980*. arXiv: 1412.6980. URL: <http://arxiv.org/abs/1412.6980>.
- Kingma, Diederik P and Max Welling (2013). "Auto-encoding variational bayes". In: *arXiv preprint arXiv:1312.6114*.
- Kohl, Nate and Peter Stone (2004). "Policy gradient reinforcement learning for fast quadrupedal locomotion". In: *Robotics and Automation, 2004. Proceedings. ICRA'04. 2004 IEEE International Conference on*. Vol. 3. IEEE, pp. 2619–2624.
- Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E Hinton (2012). "Imagenet classification with deep convolutional neural networks". In: *Advances in neural information processing systems*, pp. 1097–1105.
- Lemar  chal, Claude (2012). "Cauchy and the gradient method". In: *Doc Math Extra* 251, p. 254.
- Mnih, Volodymyr et al. (2015). "Human-level control through deep reinforcement learning". In: *Nature* 518.7540, p. 529.
- Mnih, Volodymyr et al. (2016). "Asynchronous methods for deep reinforcement learning". In: *International Conference on Machine Learning*, pp. 1928–1937.
- Monahan, George E (1982). "State of the art—a survey of partially observable Markov decision processes: theory, models, and algorithms". In: *Management Science* 28.1, pp. 1–16.
- Ng, Andrew Y et al. (2006). "Autonomous inverted helicopter flight via reinforcement learning". In: *Experimental Robotics IX*. Springer, pp. 363–372.
- OpenAI Gym. <https://gym.openai.com/envs/FrozenLake-v0/>. (Accessed on 01/25/2018).
- openai/gym on Github. <https://github.com/openai/gym/tree/master/gym>. (Accessed on 01/25/2018).
- Pan, Sinno Jialin and Qiang Yang (2010). "A survey on transfer learning". In: *IEEE Transactions on knowledge and data engineering* 22.10, pp. 1345–1359.
- Radford, Alec, Luke Metz, and Soumith Chintala (2015). "Unsupervised representation learning with deep convolutional generative adversarial networks". In: *arXiv preprint arXiv:1511.06434*.
- Schulman, John et al. (2015). "Trust region policy optimization". In: *International Conference on Machine Learning*, pp. 1889–1897.
- Shvachko, Konstantin et al. (2010). "The hadoop distributed file system". In: *Mass storage systems and technologies (MSST), 2010 IEEE 26th symposium on*. Ieee, pp. 1–10.
- Singh, Satinder et al. (2002). "Optimizing dialogue management with reinforcement learning: Experiments with the NJFun system". In: *Journal of Artificial Intelligence Research* 16, pp. 105–133.
- Strehl, Alexander L et al. (2006). "PAC model-free reinforcement learning". In: *Proceedings of the 23rd international conference on Machine learning*. ACM, pp. 881–888.
- Sutskever, Ilya et al. (2013). "On the importance of initialization and momentum in deep learning". In: *International conference on machine learning*, pp. 1139–1147.
- Tesauro, Gerald (1995). "Temporal difference learning and TD-Gammon". In: *Communications of the ACM* 38.3, pp. 58–68.
- Watkins, Christopher JCH and Peter Dayan (1992). "Q-learning". In: *Machine learning* 8.3-4, pp. 279–292.
- Zamora, Iker et al. (2016). "Extending the OpenAI Gym for robotics: a toolkit for reinforcement learning using ROS and Gazebo". In: *arXiv preprint arXiv:1608.05742*.