



CHALMERS
UNIVERSITY OF TECHNOLOGY



UNIVERSITY OF GOTHENBURG

Improving sample-efficiency of model-free reinforcement learning algorithms by learning latent space representations

a

Master's thesis in Computer science and engineering

NAME FAMILYNAME

MASTER'S THESIS 2022

Improving sample-efficiency of model-free reinforcement learning algorithms by learning latent space representations

a

NAME FAMILYNAME



UNIVERSITY OF
GOTHENBURG



CHALMERS
UNIVERSITY OF TECHNOLOGY

Department of Computer Science and Engineering
CHALMERS UNIVERSITY OF TECHNOLOGY
UNIVERSITY OF GOTHENBURG
Gothenburg, Sweden 2022

Improving sample-efficiency of model-free reinforcement learning algorithms by learning latent space representations

a

NAME FAMILYNAME

© NAME FAMILYNAME, 2022.

Supervisor: Name, Department

Advisor: Name, Company or Institute (if applicable)

Examiner: Name, Department

Master's Thesis 2022

Department of Computer Science and Engineering

Chalmers University of Technology and University of Gothenburg

SE-412 96 Gothenburg

Telephone +46 31 772 1000

Cover: Description of the picture on the cover page (if applicable)

Typeset in L^AT_EX

Gothenburg, Sweden 2022

Improving sample-efficiency of model-free reinforcement learning algorithms by learning latent space representations

a

NAME FAMILYNAME

Department of Computer Science and Engineering

Chalmers University of Technology and University of Gothenburg

Abstract

Abstract text about your project in Computer Science and Engineering.

Keywords: Computer, science, computer science, engineering, project, thesis.

Acknowledgements

Here, you can say thank you to your supervisor(s), company advisors and other people that supported you during your project.

Name Familyname, Gothenburg, May 2022

Contents

List of Figures	xi
List of Tables	xiii
1 Introduction	1
2 Background	5
2.1 Introduction to reinforcement learning	5
2.1.1 Problem setting	5
2.1.2 Bandit problems	6
2.1.3 Markov Decision Processes	6
2.1.4 Key concepts in reinforcement learning	7
2.1.4.1 Policy	7
2.1.4.2 Goal of reinforcement learning	8
2.1.4.3 Value functions	9
2.2 Classes of reinforcement learning algorithms	9
2.2.1 Policy gradients	9
2.2.1.1 Baselines	10
2.2.1.2 Off-policy gradients	11
2.2.1.2.1 Advanced policy gradient	11
2.2.2 Actor-critic algorithms	12
2.2.3 Value function methods	14
2.2.3.1 Dynamic programming	14
2.3 Deep Reinforcement Learning and DQN	14
2.3.1 Extension of DQN	15
2.3.1.1 Double Deep Q-networks: DDQN	15
2.3.1.2 Prioritized replay	15
2.3.1.3 Dueling Network	15
2.3.1.4 Multi-step learning	16
2.3.1.5 Noisy Nets	16
2.3.1.6 Integrated Agent:Rainbow	16
2.4 Problems with RL	16
2.5 general computer vision stuff	17
2.6 General latent space learning	17
2.6.1 Auto-encoder	18
2.6.2 Forward model	18

2.6.3	Inverse model	19
2.6.4	Using prior knowledge to constrain the state space	20
2.6.5	Using hybring objectives	20
2.6.6	Common neural network architectures	21
2.6.6.1	AE	21
2.6.6.2	DAE	21
2.6.6.3	VAE	21
2.6.6.4	Siamese networks	21
2.7	Model-based reinforcement learning	21
3	Related Work	23
3.1	Work on top of which we build	23
3.2	Current state-of-the-art on Atari	24
3.3	Work whose techniques we share	24
3.4	Work achieving same goals as us, but differently	24
3.5	Strongly related to our work	24
4	Methods	25
4.0.1	Enviroment and Preprocessing	25
4.0.2	Deep Auto-encoder and Model Architecture	25
4.0.3	Training the RL Agent	26
5	Results	27
5.0.1	Two Step Training	27
5.0.2	Parallel Training	27
6	Conclusion	29
6.1	Discussion	29
6.2	Conclusion	29
	Bibliography	31
A	Appendix 1	I

List of Figures

2.1	Conceptual schematic of reinforcement learning.	5
2.2	Schematic of a Markov chain.	6
2.3	Schematic of a Markov decision process.	7
2.4	Schematic of a partially observable Markov decision process.	7
2.5	Auto-encoder: learned by reconstructing the observation (one-to-one). The observation is the input and the computed state is the vector at the auto-encoder's bottleneck layer, i.e. is the output of the encoder part of the auto-encoder network. The loss is calculated between the true observation and the reconstructing observation (which is obtained by passing the observation through both the encoder and the decoder).	19
2.6	Forward model: predicting the future state from the state-action pair. The loss is computed from comparing the predicted state against the true next state (the states being the learned states). This can also be done directly by predicting the next observation and comparing against it.	19
2.7	Inverse model: predicting the action between two consecutive states. The loss is computed from comparing the predicted action between two consecutive states against the true action that was taken by the agent between those two states. (the states being the learned states).	20

List of Tables

2.1	Li Zhou ,relationship among bandit, MDP, RL, and decision theory	6
-----	--	---

1

Introduction

In computer science, reinforcement learning is the formalization of trial-and-error learning. While this is not the only legitimate interpretation of the concept, it is the most straightforward one: “trial” implies existence of an agent which observes its environment and interacts with it through its own actions. “Error” implies that the agent has a goal it tries to achieve and that it does not know how to achieve it (in the most effective manner). What it can do is take different actions and appraise them according to how closely they lead the agent toward its goal, thereby observing the quality of those actions. By repeatedly exploring the effects of various sequences of actions, the agent can find, i.e. learn, the sequence of actions which lead to its goal. Here it is important to discuss what a goal is. To formalize the process outlined above, one needs to describe it in purely mathematical terms. Thus, among other things, the goal needs to be described numerically. To achieve that, the notion of a reward function is used: it maps every state of the environment to a number which denotes its value called the reward. The state of the environment to which the highest reward is ascribed is then the goal. A more general description of the goal of reinforcement learning is to maximize the sum of rewards over time. Formalization of the entire process will be carried out later in the text, while here only the most important concepts will be outlined. One of these is the trade-off between “exploration” and “exploitation.” To learn just by trial and error implies learning from experience. This means that the agent can not know how a certain strategy fares unless it collects experiences which come by following said strategy. Thus in order to find a good strategy, usually referred to as a “policy”, the agent needs to produce various different strategies and observe their results until it finds a promising one. The process of finding different strategies and experimenting with new random behavior is called exploration. Likewise, the process of repeating a good strategy is called exploitation. Due to the curse of dimensionality,¹ in complex multi-dimensional domains it is impossible to test but a miniscule proportion of all possible strategies. Because of this, the problem of effective exploration and the trade-off between it and exploitation is a fundamental problem in reinforcement learning.

Due to its generality, reinforcement learning is studied in many different disciplines: control theory, game theory, information theory, simulation-based optimization, multi-agent systems etc.. Of these, control theory is of particular importance because it often enables clear analysis of various reinforcement learning algorithms. This foremost concerns the usage of dynamic programming which provides a basis

¹The curse of dimensionality refers to the exponential rise of possible configurations of the problem with the number of dimensions.

for a large class of reinforcement learning algorithms. Reinforcement learning is also considered to be one of the pillars of modern data-driven machine learning. In the context of machine learning, reinforcement learning can be viewed as a combination of supervised and unsupervised learning: the “trial” portion of the trial-and-error learning can be interpreted as unsupervised or as self-supervised learning because in it the agent collects its own dataset without any explicit labels to guide its way. This process is referred to as “exploration”. The dataset created by exploration is labelled by the reward function. Thus the agent can learn from “past experience” in a supervised manner. This text will introduce concepts from both control theory and machine learning which are necessary to formalize the reinforcement learning objective and to develop algorithms to achieve it. It will not concern itself with other disciplines.

Interest in reinforcement learning has grown tremendously over the past decade. It has been fueled by successes of deep machine learning in fields such as computer vision. The subsequent utilization of neural networks in reinforcement learning, dubbed deep reinforcement learning, led to impressive results such as achieving better-than-human performance on Atari games, in the game of go and in many others. Because large amounts of data are required for neural network training and thus for reinforcement learning algorithms which utilize them, most of these results are achieved in computer-simulated environments.² These recent successes were kick-started by Deep Q-Networks (DQN) algorithm which, by utilizing convolutional neural network, crucially enabled the agents to successfully learn from raw pixels. Learning from pixels is incredibly important for many practical applications, such as those in robotics where it is usually impossible to get full access to the state of the environment. The state then needs to be inferred from observations such as those from cameras.³ Due to incredible results achieved in simulated environments, reinforcement learning holds the promise of solving many incredibly important engineering problems, for example robotic manipulation and grasping. Having that said, there exists a large gap between simple simulated environments and the real world, and many improvements to the current state-of-the-art algorithms are required to bridge that gap. To explain the approach investigated in this thesis, a bit more context is needed.

An important classification of reinforcement learning algorithm is the one between model-based and model-free algorithms. As the name suggests, model-free algorithms do not form an explicit model of the environment. Instead, they function as black-box optimization algorithms, simply finding actions which maximize reward without other concerns such as predicting the states resulting from those actions. In other words, they only predict the reward of actions in given states. Model-based algorithms on the other hand learn an explicit model of the environment and use it to plan their actions. They thus learn the dynamics of the environment and use that

²Simulated environments run as fast as the computers they run on, which enables generating thousands of trials in seconds.

³Here the state refers to the underlying physical parameters of the environment: the positions and velocities of objects, the friction coefficients and so on. Observations from sensors such as cameras do not explicitly provide such information. However, since humans and animals are able to utilize such observations to achieve their goals, we know that they implicitly hold enough information about the true state of the world for successful goal completion.

knowledge to choose actions which lead the agent to states with high reward. Both classes have their benefits and their drawbacks. Since model-free algorithms do not require any knowledge of environment dynamics to operate, they are more widely applicable and usually achieve better performance. But the fact that they can not leverage environment dynamics to create plans implies a harder learning problem: they need to implicitly learn those dynamics while only being provided the reward signal. This makes them much less sample-efficient. By contrast, model-based algorithms are of course more sample-efficient. Furthermore, the plan generated from the learned model can be utilized to interpret the agent's actions which in turn leads to many further benefits such as the ability to guarantee outcomes in safety-critical operations. Unfortunately, the twin learning objective of learning the best action-choosing policy to maximize the reward over time, and the learning of the model results in fundamental training instabilities which usually results in worse final performance. In simple terms, the reason behind this is the following one: in the beginning of the learning process, both the policy and the model perform poorly. For the model to perform better, the agent needs to explore the environment and update its model. However, many parts of the environment are inaccessible to a poorly performing agent: for example, if an agent is playing a computer game, and it is not able to progress to further sections of the game, it will not be able to construct a model of that portion of the game. Thus, to explore the environment and improve its model, it needs to first learn to exploit the model and perform sufficiently well using it. Furthermore, what it learned at this stage may become obsolete as the model changes. How bad this problem is depends on the specifics of the setting, and there are many ways to ameliorate it, but in most cases the necessary trade-offs result in a lower final performance.

Given the previous discussion, the goal of the thesis may be presented: the idea is to combine the sample-efficiency of model-based approaches with the flexibility of model-free methods. Another way to describe the same is to say that we want to utilize learning signals other than the reward signal to make the model-free learning more sample-efficient. In particular, we want to learn a latent representation of the environment, i.e. to find lower-dimensional embeddings of the environment, and learn a policy in this space. To make this a concrete and manageable goal, we constrain ourselves to the problem of learning from images in particular. To be able to compare our results to those of other researchers, we will test our algorithms on the standard benchmark tasks in the field, namely Atari57 games. The potential benefits of the proposed approach are two-fold: firstly, we know that in general lower-dimensional optimization problems are easier to solve than higher-dimensional ones. Secondly, it is known that when algorithms learn with direct state access, they learn much faster and often achieve better final results. The main reason behind this is that images are much higher-dimensional than underlying states, and this is self-evident in the case of Atari games. Since inferring states from observations is not directly related to the reward, we expect that using unsupervised learning techniques will aid in feature extraction and thus make learning more sample-efficient. Furthermore, since the goal is not to learn the dynamics of the environment, but simply to find an equivalent, but lower-dimensional representation of it, we expect that this approach won't suffer from the problems faced by model-based approaches.

Of course, we are not the first to suggest such an approach, but we haven't found a systematic analysis of it accross an array of model-free algorithms. Related work will be discussed in its own chapter.

In total, our contributions will be the following:

1. A systematic analysis of utilizing autoencoders for learning lower-dimensional representations for greater sample-efficiency in combination with different reinforcement learning algorithms.
2. Ideally, finding algorithms which achieve the same or better performance with the autoencoder as without it, but doing so more efficiently.
3. Defining a clear direction for further research based on the previous two contributions.

2

Background

2.1 Introduction to reinforcement learning

2.1.1 Problem setting

In the usual engineering approach to problems, prior scientific knowledge is used to first describe the problem and then to define it mathematically. Once this is done, unknown variables are measured and solutions are calculated. This approach works if the inherent stochasticity of the environment can be controlled, i.e. if bounds of stochasticity are known the solution account for them and be designed to be robust to them. But some problems have circumstances which can not be known in advance, or which are incredibly hard to hand-engineer.

In those cases, an entirely different approach becomes the only viable one: designing a system which can produce and refine its own solution, or in other words, designing a system which, in a way, learn the solution by itself. This is the idea behind the learning-based approach: automating the process of learning. Crucially, now the world and how it operates is unknown and has to be discovered. The schematic 2.1 shows how this process is formulated in reinforcement learning. Reinforcement

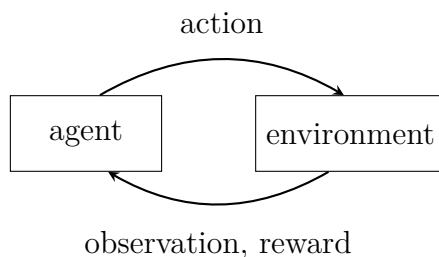


Figure 2.1: Conceptual schematic of reinforcement learning.

learning is a 2-step iterative process. The **agent**, which represents the computer program, takes **actions** in its **environment**. It then **observes** the resulting state of the environment and is also given a **reward** which is a function mapping every state of the environment to a number.

To introduce reinforcement learning more formally, we first describe the simplest possible problem to which reinforcement learning is the best solution.

2.1.2 Bandit problems

Reinforcement learning uses training information that evaluates the actions taken rather than instruct by giving correct actions. Consider this learning problem. The agent is faced with k different gambling slot machines. Each of them give random rewards under an unknown distribution. At each turn, the agents has to select one of the machines and pull its lever. The goal is to maximize the expected total reward over some number of turns. If the agent knew the distribution of rewards of each of the slot machines, it would simply choose the one with the highest expected reward in number of turns it has been given. At any time step the agent will be able to select at least one action whose estimated value is greatest. When the agent selects on of this actions it is exploiting the current knowledge of the values of the actions. If the agents keep exploiting the goal of maximizing reward over period of time will be trivial. If instead the agent selects one of the non greedy action this will enables it to improve the average expected reward over time.

When addressing the canonical problem of sequential decision making under uncertainty, the exploitation-exploration trade-off is highlighted. More specifically, as depicted in Fig.1, an agent interacts with an unknown environment in a sequential manner to obtain rewards. The ultimate goal is to maximize the rewards. For one thing, the agent takes advantage of existing knowledge of the environment. For another, the agent investigates an unfamiliar environment.

Table 2.1: Li Zhou ,relationship among bandit, MDP, RL, and decision theory

<i>Learn model outcome</i>	Mutlti armed bandits	RL
<i>Given model of stocastic outcomes</i>	Decision Theory	MDP
	Actions dont change state of the world	Actions change state of the world

2.1.3 Markov Decision Processes

The environment of k -bandit problem is static — the actions do not change the **state** of the environment. To model environments in which states change, Markov chains are used. They capture the stochastic nature of state transitions, while Markov property allows for easier mathematical analysis. The schematic of a Markov chain is shown in 2.3.

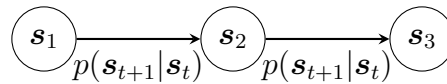


Figure 2.2: Schematic of a Markov chain.

Formally, a Markov chain \mathcal{M} is defined by its state space \mathcal{S} with discrete or continuous state $\mathbf{s} \in \mathcal{S}$ and the transition operator \mathcal{T} . The notation \mathbf{s}_t denotes the state at time t and it is a vector of real numbers. The transition operator allow for a succinct description of environment dynamics. For a transition probability $p(\mathbf{s}_{t+1}|\mathbf{s}_t)$, let $\mu_{t,i} = p(\mathbf{s}_t = i)$ and $\mathcal{T}_{i,j} = p(\mathbf{s}_{t+1} = i|\mathbf{s}_t = j)$. Then $\vec{\mu}_t$ is a vector of probabilities and $\vec{\mu}_{t+1} = \mathcal{T}\vec{\mu}_t$. Importantly, \mathcal{T} is linear.

To model the agent's actions, we simply augment the Markov chain by adding actions as priors to state transition probabilities and defining the reward function, thereby constructing a Markov decision process. It's schematic can be seen in ??.

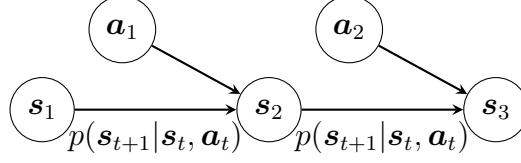


Figure 2.3: Schematic of a Markov decision process.

The Markov decision process is thus defined by the tuple $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{T}, r\}$. \mathcal{A} denotes the action space, where $\mathbf{a} \in \mathcal{A}$ is a continuous or discrete action and r is the reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. It should also be noted that now the transition operator is a tensor. Let $\mu_{t,j} = p(s_t = j)$, $\xi_{t,k} = p(a_t = k)$, $\mathcal{T}_{i,j,k} = p(s_{t+1} = i | s_t = j, a_t = k)$. Then $\mu_{t+1,i} = \sum_{j,k} \mathcal{T}_{i,j,k} \mu_{t,j} \xi_{t,k}$. Hence \mathcal{T} retains its linearity.

Finally, partial observability also needs to be accounted for. To do so, a partially observable Markov decision process (POMDP) needs to be constructed. This is done by augmenting the Markov decision process to also include the observation space \mathcal{O} , where observations $\mathbf{o} \in \mathcal{O}$ denote the discrete or continuous observations and the emission probability \mathcal{E} which describes the probability $p(\mathbf{o}_t | \mathbf{s}_t)$ of getting the observation \mathbf{o}_t when in state \mathbf{s}_t . The schematic can be seen in 2.4.

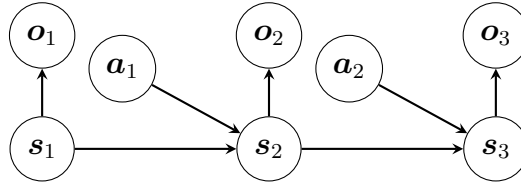


Figure 2.4: Schematic of a partially observable Markov decision process.

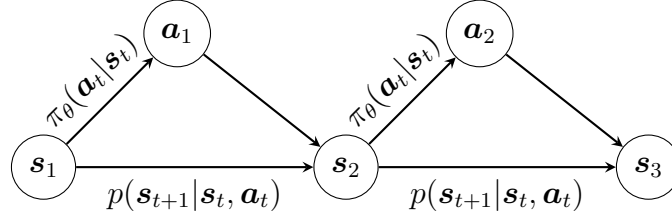
It is important to note that not all elements of POMDP are present in every problem: for example, the reward may be a deterministic function of the state and so on. In general through the text, to aid in simplifying notation, only the necessary elements will be explicitly referenced in sketches and written out in the equations (most often using just the Markov decision process).

2.1.4 Key concepts in reinforcement learning

2.1.4.1 Policy

With the problem space being formally defined, we may introduce definitions which will allow the construction of reinforcement learning algorithm. The reinforcement learning problem can be defined in finite or infinite time horizons. Different environments usually naturally fall in either category. For the agent to learn, it needs to be able to try out different actions from the same, or at least similar states. This is usually achieved by having the agent return to a set of starting states. The period between two such returns is called **an episode**. The agent selects actions based on

its **policy** π . The policy is a function which maps states to actions. The schematic showing it in the context of a Markov decision process is given in ??.



The policy is a stochastic function. The intensity of stochasticity determines the trade-off between exploration and exploitation. To emphasize that the policy depends on some parameters θ , we usually write π_θ .

2.1.4.2 Goal of reinforcement learning

For simpler notation, the finite horizon form is assumed for the following definitions. Since the environment is modelled as a Markov decision process, we can write the probability of observing a trajectory of states and actions as:

$$\underbrace{p_\theta(\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_T, \mathbf{a}_T)}_{p_\theta(\tau)} = p(\mathbf{s}_1) \prod_{t=1}^T \underbrace{\pi_\theta(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)}_{\text{Markov chain on } (\mathbf{s}, \mathbf{a})} \quad (2.1)$$

A bit more explicitly, we can a transition probability as:

$$p((\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) | (\mathbf{s}_t, \mathbf{a}_t)) = p((\mathbf{s}_{t+1} | (\mathbf{s}_t, \mathbf{a}_t)) \pi_\theta(\mathbf{a}_{t+1} | \mathbf{s}_{t+1}) \quad (2.2)$$

With this, we may now formally define the goal of reinforcement learning. It is to find policy parameters θ^* such that:

$$\theta^* = \operatorname{argmax}_{\theta} E_{\tau \sim p_\theta(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \quad (2.3)$$

$$= \operatorname{argmax}_{\theta} \sum_t^T E_{(\mathbf{s}_t, \mathbf{a}_t) \sim p_\theta(\mathbf{s}_t, \mathbf{a}_t)} [r(\mathbf{s}_t, \mathbf{a}_t)] \quad (2.4)$$

To ensure that the expected sum of rewards, also know as the **return**, is finite in the infinite horizon case, a **discount factor** $0 < \gamma < 1$ is introduced in the sum. The discount factor also play a role in modelling because often times it makes sense to value immediate rewards more. It is important to note that we are maximizing the *expected* sum of rewards. This it a the goal a smooth and differentiable function of the parameters, which means we can employ gradient descent to find the optimal parameters. This leads us to the first class of reinforcement learning algorithms: policy gradient algorithms. They will be introduced with the other classes of algorithm in the next subsection, while here additional concepts required by other classes of algorithms will be introduced here.

2.1.4.3 Value functions

Value functions are functions which map states or state-action pairs to the expected returns obtained under a fixed policy. They are a concept from dynamic programming. In fact, reinforcement learning can be interpreted as an extension of dynamic programming, as shall be done in the following subsection. Having that said, value function can be interpreted in other ways as well. The **Q-function** maps state-action pairs to the estimated sum of returns under policy π_θ :

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t] \quad (2.5)$$

thus denoting the total reward from taking \mathbf{a}_t in \mathbf{s}_t . **Value functions** map states to to the estimated sum of returns under policy π_θ :

$$V^\pi(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'} | \mathbf{s}_t)] \quad (2.6)$$

The connection between the two is the following:

$$V^\pi(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi(\mathbf{s}_t, \mathbf{a}_t)} [Q^\pi(\mathbf{s}_t, \mathbf{a}_t)] \quad (2.7)$$

And we can also write the RL objective as:

$$E_{\mathbf{s}_1 \sim p(\mathbf{s}_1)} [V^\pi(\mathbf{s}_1)] \quad (2.8)$$

2.2 Classes of reinforcement learning algorithms

2.2.1 Policy gradients

Policy gradients are derived by directly solving for the reinforcement learning objective with gradient descent with respect to the policy parameters. To do so, the reinforcement learning objective needs to be evaluated. We begin by introducing a notational shorthand:

$$\theta^* = \operatorname{argmax}_{\theta} \underbrace{E_{\tau \sim p_\theta(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right]}_{J(\theta)} \quad (2.9)$$

We estimate $J(\theta)$ by making rollouts from the policy (below i is the sample index and i, t is the t^{th} timestep in the i^{th} sample):

$$J(\theta) = E_{\tau \sim p_\theta(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \quad (2.10)$$

Simplifying the notation further, we get:

$$J(\theta) = E_{\tau \sim p_\theta(\tau)} \underbrace{[r(\tau)]}_{\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t)} = \int p_\theta(\tau) r(\tau) d\tau \quad (2.11)$$

The goal now is to compute the derivative of the estimated reinforcement learning objective:

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau \quad (2.12)$$

Since the goal of this text is just to introduce the necessary concepts and algorithms, the derivation(s) will be omitted. We encourage the interested reader to consult the literature [SB18] and CITE LEVINE'S BERKLEY LECTURES to find them. Here we will just note that it is crucial that the final expression can be estimated by sampling the agent's experience as the other quantities are not available. The resulting expression for the policy gradient is:

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right] \quad (2.13)$$

To evaluate the policy gradient we can sample:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right) \quad (2.14)$$

With the gradient we can do a step of gradient ascent and use it to form the REINFORCE algorithm, also known as “vanilla policy gradient”:

REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ by running the policy
2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

This algorithm does not work well in practice. The main reason for that is that the variance of returns is very high. However, there are a number of modification which dramatically improve its performance. Since the goal of this text is not to outline every reinforcement learning algorithm, we will introduce only the modifications which outline general trade-offs and principles in reinforcement learning algorithm design.

2.2.1.1 Baselines

The policy gradient in the REINFORCE algorithm lacks some important properties. One of them is that it should, ideally, make bad actions less likely and good actions more likely. However, if all rewards are positive, then all action's probabilities will be increased, only by different amounts. This can be changed if a **baseline** b is added to actions:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b] \quad (2.15)$$

$$b = \frac{1}{N} \sum_{i=1}^N r(\tau) \quad (2.16)$$

This addition doesn't change the gradient in expectation, i.e. it does not introduce bias, but it does change its variance. While an optimal bias can be calculated, it is rarely used in practise due to its computational cost. Using baselines is one of the key ideas in actor-critic algorithms so they will be discussed further there.

2.2.1.2 Off-policy gradients

An important property of the REINFORCE algorithm is that it is an **on-policy** algorithm. This means that new samples need to be collected for every gradient step. The reason behind this is the fact that the expectation of the gradient of the return needs to be calculated with respect to the current parameters of the policy. In other words, because the policy changes with each gradient step, old samples are effectively collected under a different policy. This means that they can not be used to calculate the expected gradient of the return with respect to the current policy — it would not produce those trajectories. In mathematical notation:

$$\nabla_{\theta} J(\theta) = \underbrace{E_{\tau \sim p_{\theta}(\tau)}}_{\text{this is the trouble!}} [\nabla_{\theta} p_{\theta}(\tau) r(\tau)] \quad (2.17)$$

If the policy is a neural network, which requires small gradient steps, the cost of generating a big number of samples for every update could make the algorithm entirely infeasible. This of course depends on the cost of generating samples, which is entirely problem dependent — policy gradient algorithms are often the best solution when the cost of generating samples is low.

However, on-policy algorithms can be turned into off-policy algorithms through **importance sampling**, which is the name given to the following mathematical identity:

$$E_{x \sim p(x)}[f(x)] = \int p(x) f(x) dx \quad (2.18)$$

$$= \int \frac{q(x)}{q(x)} p(x) f(x) dx \quad (2.19)$$

$$= \int q(x) \frac{p(x)}{q(x)} f(x) dx \quad (2.20)$$

$$= E_{x \sim p(x)} \left[\frac{p(x)}{q(x)} f(x) \right] \quad (2.21)$$

which is exact in expectation. To use importance sampling to create an off-policy policy gradient algorithm, certain approximations need to be made. Again, the details of the derivation are omitted and what follows is just the final result.

$$\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \hat{Q}_{i,t} \quad (2.22)$$

To get this equation, the factor $\frac{\pi_{\theta'}(\mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{s}_{i,t})}$ had to be simply ignored in the expression because it is impossible to calculate the state marginal probabilities. This means that the expression works only if $\pi_{\theta'}$ is not too different from π_{θ} .

2.2.1.2.1 Advanced policy gradient The basic algorithm we have outlined is essentially just a basic gradient descent method. From convex optimization, we know that it can be made much better if second order derivatives or their approximations are used. For example, conjugate gradient descent can be used. Further, there are various ways in which this optimization problem can be better conditioned. Such improvements led algorithms such as PPO and TRPO, which will not be discussed here.

2.2.2 Actor-critic algorithms

Actor-critic methods can be seen as making a different trade-off between variance and bias in policy gradient estimation. We begin with the following observation: ¹

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \underbrace{\left(\sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)}_{\hat{Q}_{i,t}: \text{“reward to go”}} \quad (2.23)$$

Simply put, in the policy gradient method a single-run Monte-Carlo (MC) is used to estimate the return. This causes high variance, while incurring no bias. Another option is to try to estimate the full expectation $\hat{Q}_{i,t} \approx \sum_{t'=t}^T E_{\pi_{\theta}} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$. Since the estimate won't be perfect, it will introduce bias. Of course, using multiple runs from the same state-action pair would reduce variance, but this is sometimes impossible to procure and is certainly more costly. However, if our estimator of “reward to go” can generalize between states, we will be able to get good estimates regardless.

Like the policy, the return estimator will have to be learned. In this approach, the policy is also called the **actor** and the return estimator is called the **critic**. We proceed by discussing how the critic can be constructed. If we had the correct Q-function (i.e. not the estimate, but the actual values), we could improve the policy gradient estimate by using it both to estimate the return and as a baseline:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) (Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) - b) \quad (2.24)$$

$$b_t = \frac{1}{N} \sum_i Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \quad (2.25)$$

However, having a baseline that depends on actions leads to bias. Thus we employ a baseline dependent on the state:

$$V(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_{\theta}(\mathbf{s}_t, \mathbf{a}_t)} [Q(\mathbf{s}_t, \mathbf{a}_t)] \quad (2.26)$$

Since the value function 2.6 tells us the expected return of the average action, we can calculate how much better a certain action is by subtracting its Q-value 2.5 for the value function. The result is called the **advantage function**:

$$A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) - V^{\pi}(\mathbf{s}_t) \quad (2.27)$$

Thus we can fit either the Q-function, the value function or the advantage function. Of these, it is best to learn the value function because there are less states than state-action pairs. We then calculate the advantage function in the following way:

$$A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^{\pi}(\mathbf{s}_{t+1}) - V^{\pi}(\mathbf{s}_t) \quad (2.28)$$

¹In this equation, the summation of rewards is done from time t to T because actions and states prior to that time do not affect the return from that time onward. This leveraging of causality reduces the variance of the estimate.

The value function can be estimated through samples

$$V^\pi(\mathbf{s}_t) \approx \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \quad (2.29)$$

After collecting many such samples

$$\left\{ \left(\mathbf{s}_{i,t}, \underbrace{\sum_{t'=t}^T r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})}_{y_{i,t}} \right) \right\} \quad (2.30)$$

we can fit the value function through supervised regression with the loss being:

$$\mathcal{L}(\phi) = \frac{1}{2} \sum_i \|\hat{V}_\phi^\pi(\mathbf{s}_i) - y_i\|^2 \quad (2.31)$$

However, this process can be sped up with bootstrapped estimates:

$$y_{i,t} = \sum_{t'=t}^T E_{\pi_\theta} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{i,t}] + V^\pi(\mathbf{s}_{i,t+1}) \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_\phi^\pi(\mathbf{s}_{i,t+1}) \quad (2.32)$$

This will further reduce variance, but again increase bias.

Fortunately, we can tune the trade-off between bias and variance. In the Monte Carlo estimate, the entire trajectory was used to estimate the return. In the bootstrap estimate, only a single step in the future was used along with the estimate. Instead, a **n-step** return estimator can be used:

$$\hat{A}_n^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{t+n} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) - \hat{V}_\theta^\pi(\mathbf{s}_t) + \gamma^n \hat{V}_\theta^\pi(\mathbf{s}_{t+n}) \quad (2.33)$$

In most cases the ideal trade-off for n lies somewhere between 1 and ∞ (the MC estimate). Finally, an average of all n-step return estimators can be used. This is called the generalized advantage estimator (GAE):

$$\hat{A}_{GAE}^\pi(\mathbf{s}_t, \mathbf{a}_t) = \sum_{n=1}^{\infty} (\gamma\lambda)^{n-1} \left(r(\mathbf{s}_{t+n}, \mathbf{a}_{t+n}) - \hat{V}_\theta^\pi(\mathbf{s}_{t+n}) + \gamma \hat{V}_\theta^\pi(\mathbf{s}_{t+n-1}) - \hat{V}_\theta^\pi(\mathbf{s}_t) \right) \quad (2.34)$$

where the factor λ controls the weight of future values.

Combining this into an iterative algorithm, and fixing the issues of naive implementations results in the following algorithm:

Actor-critic algorithm template

1. take action $\mathbf{a} \sim \pi_\theta(\mathbf{a} | \mathbf{s})$, get $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$, store in \mathcal{R} (replay buffer)
2. sample a batch $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$ from buffer \mathcal{R}
3. update \hat{Q}_θ^π using target $y_i = r_i + \gamma \hat{Q}_\theta^\pi(\mathbf{s}'_i, \mathbf{a}'_i) \forall \mathbf{s}_i, \mathbf{a}_i$
4. $\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(\mathbf{a}_i | \mathbf{s}_i) \hat{Q}_\theta^\pi(\mathbf{s}_i, \mathbf{a}_i)$, where $\mathbf{a}_i^\pi \sim \pi_\theta(\mathbf{a} | \mathbf{s}_i)$
5. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

2.2.3 Value function methods

Value function methods use only the critic from actor-critic algorithms. Suppose that the advantage function $A^\pi(\mathbf{s}_t, \mathbf{a}_t)$ is known. It tells us how much better the action \mathbf{a}_t is than the average action according to the policy π . Thus, if we knew the advantage function, we could construct a deterministic **greedy policy**:

$$\pi_{\text{greedy}}(\mathbf{s}_t|\mathbf{a}_t) = \begin{cases} 1 & \text{if } \mathbf{a}_t = \operatorname{argmax}_{\mathbf{a}_t} A^\pi(\mathbf{s}_t, \mathbf{a}_t) \\ 0 & \text{otherwise} \end{cases} \quad (2.35)$$

which would yield the highest expected return. In other words, if we knew the advantage function, the policy would be reduced to the argmax operation.

2.2.3.1 Dynamic programming

Dynamic programming refers to a collection of algorithms that can be used to compute optimal policies given a perfect model of the environment as an MDP. They are of limited utility in reinforcement learning due to the perfect model requirement and their great computational expense, but are important theoretically — they provide an essential foundation for understanding the other methods. Usually a finite MDP is assumed. DP can be applied to continuous problems as well, but exact solution exist only in special cases.

TODO:

1. Introduce value iteration in a tabular setting
2. introduce fitted value iteration
3. show that q-learning is off-policy
4. show that q-learning does not converge when using non-linear function approximation (introduce the bellman operator as well)

2.3 Deep Reinforcement Learning and DQN

(paraphrase)After deep learning has shown remarkable results in learning from raw pixel data in computer vision it's application has been adopted to solve many classical learning problems. The goal of Deep Reinforcement Learning is to connect reinforcement learning algorithm to deep neural network which operates directly on RGB images and efficiently process training data by using stochastic gradient updates[2].

The use of label data and the assumption the distribution of data to be identical and independent through out training process in deep learning makes it complex to use it directly into reinforcement learning algorithms which must be able to learn from sparse ,noisy and delayed reward and highly correlated data. To address this issues and successfully apply deep learning to reinforcement learning [2] used experience replay mechanism[15] which randomly samples previous transitions, and thereby smooths the training distribution over many past behaviors. Convolutional neural network was used to learn successful control policies from raw video data in

complex reinforcement learning environment. The network is trained with a variant of Q-learning algorithm, with stochastic gradient descent to update the weights.

This architecture was based on Tesauro's TD-Gammon architecture [17] which updates the parameters of the network that estimates the value function directly from on-policy samples of experience. Similar to this approach the online network in DQN utilizes a technique known as experience replay [18] where the agent's experiences at each time-step, $\tau_t = (s_t, a_t, r_t, s_{t+1})$ is stored in a data set $\mathcal{D} = e_1, \dots, e_N$ pooled over many episodes into a replay memory. By drawing random samples from this pool of stored experiences the Q-learning is updated. After performing experience replay, the agent selects and executes an action according to an ϵ -greedy policy. The use of experience replay and target networks enables relatively stable learning of Q values, and led to super human performance on several Atari games.

The advantage of using Deep Q-learning over Q-learning includes allowing to have greater sample efficiency, reduced variance by randomising the sample bias and avoiding being stuck in local minimum. Drawback DQNs are it only handle discrete, low-dimensional action space.

2.3.1 Extension of DQN

Several Studies were performed to increase the performance of DQNs (Paraphrase).

2.3.1.1 Double Deep Q-networks: DDQN

DQN suffer from overestimation bias due to the maximization step in optimisation function in Q-learning. Q-learning can overestimate actions that have been tried often and the estimations can be higher than any realistic optimistic estimate. Double Q-learning [19], addresses this overestimation by decoupling, in the maximization performed for the bootstrap target, the selection of the action from its evaluation. Double Q-learning stores two Q-functions, The average of the two Q values for each action and then performed ϵ -greedy exploration with the resulting average Q values. It was successfully combined with DQN to reduce overestimations.

TODOS: DDQN target equation

2.3.1.2 Prioritized replay

The main use of replay buffer is to sample transitions with maximum probability. Both DQN and DDQN samples experiences uniformly. Prioritized replay [20] samples transitions using the maximum priority, providing a bias towards recent transitions and stochastic transitions even when there is little left to learn about them.

2.3.1.3 Dueling Network

Dueling Network was designed for value based learning, this architecture separates the representation of state-value and state-dependent action advantages without supervision [6]. It consists of two streams that represent the value and advantage functions, while sharing a common convolutional feature learning module. This network

has a single Q-learning network with two streams that replace DQN architecture[3].

$$Q(s, a; \theta, \alpha, \beta) = V(s, \theta, \beta) + A(s, a; \theta, \alpha) \quad (2.36)$$

if needed **TODO ADD EQUATION:factorization of action values**

2.3.1.4 Multi-step learning

Previously stated extension of DQN have indicated that the use deep learning has enhanced the learning capability of Q-learning. The performance of Q-learning is still limited by greedy action selection after accumulating a single reward. An alternative approach was multi-step targets:

$$y_{j,t} = \sum_{t'=t}^{t+N-1} \gamma^{t-t'} r_{j,t'} + \gamma^N \max_{\mathbf{a}_{j,t+N}} Q_{\phi'}(\mathbf{s}_{j,t+N}, \mathbf{a}_{j,t+N}) \quad (2.37)$$

A multi-step variant of DQN is then defined by minimizing the alternative loss[16],

$$R_t^{(n)} + \gamma_t^{(n)} \max_{a'} q_{\theta}^-(S_{t+n}, a') - q_{\theta}(S_t, A_t) \quad (2.38)$$

2.3.1.5 Noisy Nets

The one limitation of ϵ -greedy policy is many actions must be executed to collect the first reward. Noisy Nets proposed a noisy linear layer that combines a deterministic and noisy stream. Depending on the learning rate the network ignores to learn the noisy stream.

2.3.1.6 Integrated Agent:Rainbow

In the Rainbow architecture [16] tries combine the above six method stated above all together. Distributional loss (3) was replaced with a multi-step variant. The target distribution was constructed by contracting the value distribution in S_{t+n} according to the cumulative discount, and shifting it by the truncated n-step discounted return. multi-step distributional loss with double Q-learning by using the greedy action in S_{t+n} selected according to the online network as the bootstrap action $a \cdot t + n$, and evaluating such action using the target network.

TODO ADD EQUATION:target distribution

2.4 Problems with RL

ex. drl that matters paper

2.5 general computer vision stuff

2.6 General latent space learning

TODO: throw in citations (ex. pilco, embed2control etc) As discussed in the introduction, learning control from images is very desirable. Images, and observations in general, only implicitly provide information about the underlying state. Finding a good policy from observations, especially images, is much more difficult than finding a policy with direct state access because the state first needs to be inferred from those observations. Reinforcement learning algorithms can by themselves implicitly extract the relevant information from observations, but this at best results in much less sample-efficient training and at worst results in complete failure. Often a problem which a reinforcement learning algorithm can solve with direct state access, can not achieve any progress when provided only image observations.

It is clear from the previous section that amazing results were achieved in the field of computer vision. However, to leverage these results for the purposes of reinforcement learning, the methods in question need to be applied for state estimation. This is a drastically different problem than for example image classification. The key difference is that now dynamics need to be inferred. While neural network architectures like the convolutional neural network are able to achieve great successes in timeless problems, neural architectures like LSTMs aimed at learning from sequential data comparatively perform much worse. Results in problems such as video prediction or action classification leave much to be desired.

Having that said, the learning signal generated from for example image reconstruction loss is substantially stronger than the reward signal, especially in settings with sparse rewards where it is not present most of the time. Thus it stands to reason that somehow leveraging the learning signal from some computer vision method should aid the reinforcement learning process. One way to do this is to explicitly use such methods to learn a function which maps from observations to states and then use reinforcement learning methods these learned state representations. This approach is explored in this section, mainly with the help of the [Les+18] overview paper. Here we discuss state representation learning for control in general as this will allow for a broader contextualization of our own work. Importantly, this does not concern learning a model which can be used to achieve control through planning, although there are of course similarities between these approaches.

In general, representation learning algorithms are designed to learn abstract features that characterize data. In the simplest forms they include methods such as k nearest neighbors. In state representation learning (SRL) the learned features are of low dimension, evolve through time and are depended on actions of an agent. The last point is particularly important because in reinforcement learning, features that do not influence the agent and that can not be influenced by the agent are not relevant for the problem of optimally controlling the agent. Also, simply reducing the dimensionality of the input to a reinforcement learning agent results in a computationally easier learning problem, which can make a difference between the solution being feasible or infeasible. Ideally, state representation learning should be done in an without explicit supervision as it can then be done in tandem with the likewise

unsupervised reinforcement learning.

While we assume that state-transitions have the Markov property, partial observability denies the possibility of having a one-to-one correspondence between each observation and state — an object whose position is required may be occluded by another. Thus prior observations have affect the mapping to the current state. Images in particular also do not encode kinematic or dynamic information: to get that crucial information a sequence of images is required. Hence we define the SRL task as learning a representation $\tilde{\mathbf{s}}_t \in \tilde{\mathcal{S}}$ of dimension K with characteristics similar to those of true states $\mathbf{s}_t \in \mathcal{S}$. In particular, the representation is a mapping of the history of observation to the current state: $\tilde{\mathbf{s}}_t = \phi(\mathbf{o}_{1:t})$. Actions $\mathbf{a}_{1:t}$ and rewards $r_{1:t}$ can also be added to the parameters of ϕ . This can help in extracting only the information relevant for the agent and its task. Often the representation is learned by using the reconstruction loss; $\hat{\mathbf{o}}_t$ denotes the reconstruction of \mathbf{o}_t .

In the context of reinforcement learning, state representations should ideally have the following properties:

- have the Markov property
- be able to represent the current state well enough for policy improvement
- be able to generalize to unseen states with similar features
- be low dimensional

With this we can introduce 4 different strategies for learning latent space models: the auto-encoder, the forward model, the inverse model and the model with prior. In the figures below, the white nodes are inputs and the gray nodes are outputs. The dashed rectangles are fitted around variables with which the loss is calculated.

2.6.1 Auto-encoder

The idea behind the auto-encoder is to just learn a lower-dimensional embedding of the observation space. This should make the learning problem easier due to the dimensionality reduction. The auto-encoder may be trained to denoise the observations by passing an observation with artificially added noise to the encoder, but then calculating the reconstruction loss on the image without the added noise. Formally this can be written as

$$\mathbf{s}_t = \phi(\mathbf{o}_t; \theta_\phi) \quad (2.39)$$

$$\hat{\mathbf{o}}_t = \phi^{-1}(\mathbf{s}_t; \theta_{\phi^{-1}}) \quad (2.40)$$

where θ_ϕ and $\theta_{\phi^{-1}}$ are the parameters learned for the encoder and decoder respectively.

2.6.2 Forward model

The auto-encoder does not encode dynamic information. Since that information is necessary for control, usually a few consecutive observations (or their embeddings) are stacked and passed to the reinforcement learning algorithm. This way the information about the dynamics is implicitly provided. While doing so works, it could be made more efficient by embedding the dynamic information as well. One way to achieve this is to trained a model to predict future state representations. A model

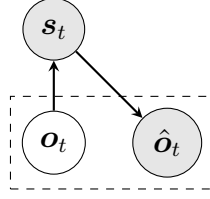


Figure 2.5: Auto-encoder: learned by reconstructing the observation (one-to-one). The observation is the input and the computed state is the vector at the auto-encoder’s bottleneck layer, i.e. is the output of the encoder part of the auto-encoder network. The loss is calculated between the true observation and the reconstructing observation (which is obtained by passing the observation through both the encoder and the decoder).

can also be observations directly, of course provided that the network in question has a bottleneck layer from which the learned representations can be extracted. Since learning on sequential information is difficult and would also benefit from lowering the dimensionality, learning a forward model can be done in two steps: first, learning an auto-encoder to embed individual frames and then learning a predictive model in the embedded space. In the schematic we show the case where predictions are learned from embeddings because it is the structurally more complex scheme. Formally, we have

$$\hat{\mathbf{s}}_{t+1} = f(\tilde{\mathbf{s}}_t, \mathbf{a}_t; \theta_{\text{forward}}) \quad (2.41)$$

The forward model can be constrained to have linear transition between $\tilde{\mathbf{s}}_t$ and $\tilde{\mathbf{s}}_{t+1}$, thereby imposing simple linear dynamics in the learned state space. Depending on the problem, if this is done well enough, learning a control law can be avoided and instead schemes like model-predictive control can be employed.

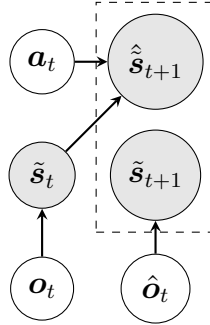


Figure 2.6: Forward model: predicting the future state from the state-action pair. The loss is computed from comparing the predicted state against the true next state (the states being the learned states). This can also be done directly by predicting the next observation and comparing against it.

2.6.3 Inverse model

The introducing predictions solves the problem of not embedding the dynamic information. However, not all information in the observation is relevant for control.

Consider a computer game where images feature decorative backgrounds — those decorations are irrelevant for playing the game well. If the reconstruction loss is computed from entire observation, that information is also carried over into the embedded space. However, if the model is trained to predict actions, it is only incentivised to use information which the agent can affect. Thus, due to less information being required, the inverse model should produce a more compact embedding. Formally, we can write this as:

$$\hat{\mathbf{a}}_t = g(\tilde{\mathbf{s}}_t, \tilde{\mathbf{s}}_{t+1}; \theta_{\text{inverse}}) \quad (2.42)$$

If the inverse model is neural network, we can recover the embedding by discarding the last few layers and use their outputs to produce the embeddings.

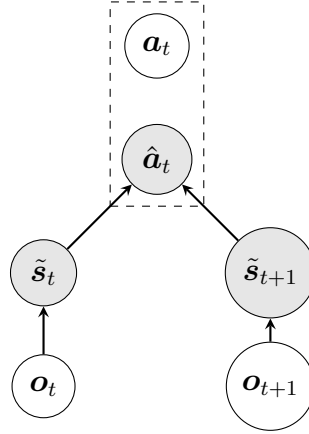


Figure 2.7: Inverse model: predicting the action between two consecutive states. The loss is computed from comparing the predicted action between two consecutive states against the true action that was taken by the agent between those two states. (the states being the learned states).

2.6.4 Using prior knowledge to constrain the state space

Of course, not everything needs to be learned in every problem. While in general hand-engineered features are worse than learned ones, there are other ways to provide prior knowledge to the learning system. For example, convolutional neural networks by their architecture encode the fact that nearby pixels are related. In the SRL context we already mention the possibility of constraining the model to linear transitions, but there are other available techniques like for example constraining temporal continuity or the principle of causality. Furthermore, priors can be defined as additional objectives or loss functions. For example, additional loss can be provided if embeddings from consecutive observations are drastically different. This is called the slowness principle.

2.6.5 Using hybrid objectives

The approaches outlined thus far can be combined into hybrid approaches. TODO: throw a reference or two from the overview paper you're going over, for example embed2control.

2.6.6 Common neural network architectures

2.6.6.1 AE

Deterministic auto-encoder

2.6.6.2 DAE

Denoising auto-encoder

2.6.6.3 VAE

Variational auto-encoder

2.6.6.4 Siamese networks

Networks that share parameters.

2.7 Model-based reinforcement learning

TODO Introduce just the idea for the sole purpose of showing why we aren't doing model-based reinforcement learning, but instead opting for model-free with state representation learning.

3

Related Work

As said in the introduction, the goal of the thesis is to use state representation learning to increase the efficiency and final results of model-free reinforcement learning. We are now ready to discuss the specifics of our approach. Firstly, we limit ourselves to image observations and discrete action spaces. In particular, we limit ourselves to Atari57 games as they are common benchmarks in the field for discrete action spaces. As shall be seen in the following text, a lot of recent work in state-representation learning for model-free reinforcement learning has been done in robotics problems with continuous action spaces. Our goal will mostly be to transfer these successes to the discrete action space setting, namely to DQN and Rainbow algorithms as the basis, instead of SAC and DDPG. Importantly, since we are concerned with finding ways to make reinforcement learning more sample-efficient, we will be using only off-policy algorithms. Secondly, we are particularly interested in the problem of simultaneous training of the state representations and the policy. The reason for this is that two-step training is often not available. This state of affairs is the natural setting for problems where reinforcement learning is a good solution: the problems where exploration is necessary due to either the high complexity of the dynamics or unanticipatable events. Parallel training of the state representations and the policy necessitates instability in policy training due to the fact the state estimations change even for same observations as the state representation are learned. Hence, related work that focuses on solving or at least ameliorating this issue is of particular importance to our work. Finally, we want our method to be robust not just in the sense that it works across a wide array of problems, but in the sense that it can be easily added to a variety of reinforcement learning algorithms and yield a positive result. In other words, it should function as a module which can be easily added to new algorithms. Furthermore, it should work well with other improvements as those suggested in some of the following related work. For clarity, we divide our discussion of related work in three categories: prior work on top of which we build, work which utilizes some of the same techniques we employ for different goal, or which achieves the same goals in different, but related ways and work which is strongly related to our own from which derived inspiration to our approach.

3.1 Work on top of which we build

1. algs we use: [Hes+18], [Haa+18]
2. (us, historic) deep Auto-Encoder Neural Networks in Reinforcement Learning [LR10]

3.2 Current state-of-the-art on Atari

1. (discuss some state of the art); [Bad+20]

3.3 Work whose techniques we share

1. (not us, but things that can be done with image loss): [She+16], [Jad+16] is the aux loss paper, also uses rewards for aux goals. point here is you can use other losses and make them intrinsic rewards.
2. (us as inspiration) [Pat+17] - here the importance of embedding only observations relevant for the agent are discussed. they used this to improve exploration though.

3.4 Work achieving same goals as us, but differently

1. (not us but same high-level idea): data-augmentation like [KYF20a], [Yar+21], [Las+20], [KYF20b], contrastive loss like [LSA20], [Sha+21], [Zha+20]
TODO First explain [Las+20] which just throws image augmentation over sampled observations. In [KYF20a], the authors go further by using multiple augmentations over the same image to regularize the Q-function. The expected return calculated for each of the augmented versions of an observation. Then the average of those is calculated and passed as the predicted expected return. This greatly stabilizes Q-learning and allows it to learn with much greater sample-efficiency. Since what is developed is a regularization technique, it can be seamlessly applied to various reinforcement learning algorithms. In [Yar+21], the authors further improve [KYF20a] by introducing several changes to their overall algorithm.
In [LSA20], contrastive loss is used instead.

3.5 Strongly related to our work

- 1.
2. (us) yarats improving sample efficiency [Yar+19], [All+21]

4

Methods

In this section, we will explain the architecture of our auto-encoder and reinforcement learning algorithm. This includes a description of the environment and preprocessing in Section 3.0.1., the collection of training data for the auto-encoder and model architecture in section 3.0.2. and the training of the RL agent in Section 3.0.3.

4.0.1 Enviroment and Preprocessing

[!ADD MORE HERE]

We perform a comprehensive evaluation of our proposed method on the Arcade Learning Environment (Bellemare et al., 2013), which is composed of 57 Atari games. The challenge is to deploy a single algorithm and architecture, with a fixed set of hyper-parameters, to learn to play all the games given embedded latent space representation of the environment from auto encoder and game rewards. This environment is very demanding because it is both comprised of a large number of highly diverse games and the observations are high-dimensional.

Working with raw Atari frames, which are 210 x 160 pixel pictures with a 128 color palette, is computationally expensive, therefore we do a basic preprocessing step to reduce the input dimensionality. The raw frames are preprocessed by down sampling to a 110 x 84 picture and transforming their RGB representation to gray-scale. Cropping an 84 x 84 rectangle of the image that nearly captures the playing area yields the final input representation to encoder part of the auto encoder.

4.0.2 Deep Auto-encoder and Model Architecture

[!ADD MORE HERE AFTER NAP]

The first step in the process is to collect data to train the auto-encoder. we run a data collection module to generate the 100000 frame for each stated under the result section. The raw images are transformed to tensors and then trained a variational autoencoder with the objective of re-constructing the original image fed to the network. The auto-encoder was trained for maximum 100 epochs. When reconstructing an image with a network bottleneck, the encoder is forced to compress the original image to a smaller dimensional vector in the latent space.

[By compressing the raw pixels environment to smaller dimensional vector in the latent space we aim to improve the training time it takes for the integrated Rl agent developed by [16] and the shift in the latent space representation and its impact

on the RL agent learning performance. We show this in more detail in the following sections.]

4.0.3 Training the RL Agent

[!ADD MORE HERE]

In this paper we integrate auto encoder with integrated agent called Rainbow[12], the selection of this architecture was based on its ability to outperform all the previous architectures. Our main focus was to experiment with the latent space representation of the environment. To address this we set up two experiment architectures: one two step training and parallel training.

Two step Training: First, we train the auto encoder with the preprocessed data and the weights of the encoder are saved in file for training the agent. In this set up the auto encoder is not updated such that we have a static representation of the environment.

Second, we train the integrated agent with the results from the auto encoder.

Parallel Training: In this set up the agent is trained using the dynamic state representation from the encoder as we train both the auto encoder and the integrated agent in parallel. This introduces stochasticity of the environment to the training and in real life control system we believe that this will set a new paradigm on the use of RL.

5

Results

state the Experiments

5.0.1 Two Step Training

5.0.2 Parallel Training

6

Conclusion

You may consider to instead divide this chapter into discussion of the results and a summary.

6.1 Discussion

6.2 Conclusion

Bibliography

- [LR10] Sascha Lange and Martin Riedmiller. “Deep auto-encoder neural networks in reinforcement learning”. In: *The 2010 international joint conference on neural networks (IJCNN)*. IEEE. 2010, pp. 1–8.
- [Jad+16] Max Jaderberg et al. “Reinforcement learning with unsupervised auxiliary tasks”. In: *arXiv preprint arXiv:1611.05397* (2016).
- [She+16] Evan Shelhamer et al. “Loss is its own reward: Self-supervision for reinforcement learning”. In: *arXiv preprint arXiv:1612.07307* (2016).
- [Pat+17] Deepak Pathak et al. “Curiosity-driven exploration by self-supervised prediction”. In: *International conference on machine learning*. PMLR. 2017, pp. 2778–2787.
- [Haa+18] Tuomas Haarnoja et al. “Soft actor-critic algorithms and applications”. In: *arXiv preprint arXiv:1812.05905* (2018).
- [Hes+18] Matteo Hessel et al. “Rainbow: Combining improvements in deep reinforcement learning”. In: *Thirty-second AAAI conference on artificial intelligence*. 2018.
- [Les+18] Timothée Lesort et al. “State representation learning for control: An overview”. In: *Neural Networks* 108 (2018), pp. 379–392.
- [SB18] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [Yar+19] Denis Yarats et al. “Improving sample efficiency in model-free reinforcement learning from images”. In: *arXiv preprint arXiv:1910.01741* (2019).
- [Bad+20] Adrià Puigdomènech Badia et al. “Agent57: Outperforming the atari human benchmark”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 507–517.
- [KYF20a] Ilya Kostrikov, Denis Yarats, and Rob Fergus. “Image augmentation is all you need: Regularizing deep reinforcement learning from pixels”. In: *arXiv preprint arXiv:2004.13649* (2020).
- [KYF20b] Ilya Kostrikov, Denis Yarats, and Rob Fergus. “Image augmentation is all you need: Regularizing deep reinforcement learning from pixels”. In: *arXiv preprint arXiv:2004.13649* (2020).
- [LSA20] Michael Laskin, Aravind Srinivas, and Pieter Abbeel. “Curl: Contrastive unsupervised representations for reinforcement learning”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 5639–5650.
- [Las+20] Misha Laskin et al. “Reinforcement learning with augmented data”. In: *Advances in Neural Information Processing Systems* 33 (2020), pp. 19884–19895.

- [Zha+20] Amy Zhang et al. “Learning invariant representations for reinforcement learning without reconstruction”. In: *arXiv preprint arXiv:2006.10742* (2020).
- [All+21] Arthur Allshire et al. “Laser: Learning a latent action space for efficient reinforcement learning”. In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. 2021, pp. 6650–6656.
- [Sha+21] Wenling Shang et al. “Reinforcement learning with latent flow”. In: *Advances in Neural Information Processing Systems* 34 (2021).
- [Yar+21] Denis Yarats et al. “Mastering visual continuous control: Improved data-augmented reinforcement learning”. In: *arXiv preprint arXiv:2107.09645* (2021).

A

Appendix 1

(UNFINISHED)

//UPDATE THE FORMAT LATER

1. Agent: It is an assumed entity which performs actions in an environment to gain some reward.
2. Environment (e): A scenario that an agent has to face. anything the agent cannot change arbitrarily is considered to be part of the environment.
3. Reward (R): An immediate return given to an agent when he or she performs specific action or task.
4. State (s): State refers to the current situation returned by the environment.
5. Policy (π): It is a strategy which applies by the agent to decide the next action based on the current state.
6. Value (V): It is expected long-term return with discount, as compared to the short-term reward.
7. Value Function: It specifies the value of a state that is the total amount of reward. It is an agent which should be expected beginning from that state.
8. Model of the environment: This mimics the behavior of the environment. It helps you to make inferences to be made and also determine how the environment will behave.
9. Model based methods: It is a method for solving reinforcement learning problems which use model-based methods.
10. Q value or action value (Q): Q value is quite similar to value. The only difference between the two is that it takes an additional parameter as a current action.