

Deep Q-Learning: Theoretical Insights from an Asymptotic Analysis

Arunselvan Ramaswamy*and Eyke Hüllermeier†

April 13, 2021

Abstract

Deep Q-Learning is an important reinforcement learning algorithm, which involves training a deep neural network, called Deep Q-Network (DQN), to approximate the well-known Q-function. Although wildly successful under laboratory conditions, serious gaps between theory and practice as well as a lack of formal guarantees prevent its use in the real world. Adopting a dynamical systems perspective, we provide a theoretical analysis of a popular version of Deep Q-Learning under realistic and verifiable assumptions. More specifically, we prove an important result on the convergence of the algorithm, characterizing the asymptotic behavior of the learning process. Our result sheds light on hitherto unexplained properties of the algorithm and helps understand empirical observations, such as performance inconsistencies even after training. Unlike previous theories, our analysis accommodates state Markov processes with multiple stationary distributions. In spite of the focus on Deep Q-Learning, we believe that our theory may be applied to understand other deep learning algorithms

1 Introduction

Reinforcement Learning (RL) is an important branch of machine learning, which has received increasing attention in the recent past. Roughly speaking, it considers an autonomous agent interacting with a dynamic environment, and seeks to learn a policy (prescribing actions depending on the current state of the environment) maximizing the agent’s welfare in the course of time. A popular variant of RL, called *Deep Reinforcement Learning* (DeepRL), combines the fundamental principles of RL with the power of deep learning. DeepRL has exhibited tremendous empirical success in recent years in wide ranging fields, from games [18] to self-driving cars [14].

In this paper, we focus on the popular DeepRL algorithm *Deep Q-Learning*, which was introduced in [15] and shown to achieve superhuman performance in playing ATARI video games. Q-learning is a specific approach to RL, which

*A. Ramaswamy is with the Heinz-Nixdorf Institute and the Department of Computer Science, Paderborn University, 33098 Paderborn, Germany (e-mail: arunr@mail.upb.de).

†E. Hüllermeier is with the Institute of Informatics at the University of Munich, 80538 Munich, Germany (e-mail: eyke@ifi.lmu.de).

focuses on learning the so-called Q-function to evaluate state-action pairs. In Deep Q-Learning, this function is represented by a deep neural network, called the Deep Q-Network (DQN), and learning the optimal Q-function is accomplished by minimizing the squared Bellman loss (error). DQN training typically involves repeated interactions with a simulator, or the use of historical data. In spite of its undoubted potential, Deep Q-Learning is still lacking a solid theoretical foundation. This also explains, at least partly, its slow adoption for real-world applications, although it generally performs well in a laboratory setting. The lack of a comprehensive understanding of the training process also hampers the explanation of empirical findings, such as suboptimal performance even when training is deemed sufficient.

First theoretical results include sufficient conditions for convergence of Deep Q-Learning, provided the DQN uses rectified linear units as activation functions [21]. The analysis requires strict conditions on the Bellman operator and the distribution of state Markov process. In [23], a non-asymptotic finite sample analysis of Deep Q-Learning with linear function approximation (instead of deep neural network (DNN) approximation) is presented. While studies like these focus on sufficient conditions for convergence, the focus in [1] is on characterizing conditions under which Deep Q-Learning is divergent. Although understanding divergence is of paramount importance, assuming linearity of the function approximator reduces the applicability of such results in real-world scenarios. This is because, in practice, deep neural networks, which are non-linear functions, are used as function approximators. There are many recent theoretical results that are based on the linearity of the function approximator, see for e.g., [20], [9] and [8]. In [7], the topic of efficient exploration in policy optimization is explored from a theoretical perspective. While these preliminary results are important and interesting, they do not immediately apply to Deep Q-Learning *as implemented in practice*, due to unrealistic simplifications and restrictive assumptions.

Our contributions. The performance of Deep Q-Learning strongly depends on the training procedure. Empirically, it has been observed that performance is great in some test scenarios and poor in others. The hitherto available theory does not explain this phenomenon, nor does it account for other empirical observations of similar kind. The main contribution of this paper is a comprehensive analysis of Deep Q-Learning that provides such explanations — under assumptions that are practical and verifiable.

We show that the squared Bellman loss is minimized over the set of state-action pairs, distributed in accordance with a measure obtained as a limit of a natural measure process associated with the training procedure. We also show that this limiting measure is stationary with respect to the state Markov process. Further, its empirical estimate can be used to retrain and boost performance. As stated earlier, the limiting measure is strongly shaped by the training process. It is worth mentioning that, unlike previous literature, our analysis allows for multiple stationary distributions of the state Markov process.

The most popular implementation of Deep Q-Learning involves the use of a target network. The use of such a network is shown to improve learning stability. However, it can be shown that the convergence properties of Deep Q-Learning does not change with the use a target network. Since, we focus on

convergence in this paper, and not stability, we do not consider implementations with target networks. More importantly, it has recently been shown in [12] that Deep Q-Learning that uses the “mellowmax” operator, instead of the usual “max” operator eliminates the need for target networks. They show superior performance as compared to traditional Deep Q-Learning with target networks, in many benchmark scenarios. Although, we do not explicitly consider the algorithm described [12], through appropriate modifications of the loss function our analysis can be extended to encompass this scenario as well.

Another popular implementation involves the use of a buffer memory called the *experience replay*. It stores past experiences for relearning purposes. The main analysis presented in Sections 3 and 4 do not account for the use of an experience replay. However, in Section 6, we discuss the steps involved in extending our analysis to account for this. We show that experience replay affects the quality of performance by shaping the limiting distribution. Additionally, it may aid in stabilizing the DQN training.

For our analyses, we utilize tools from the fields of stochastic approximation algorithms (SA) [6, 13], stochastic processes [10], measure theory [4], and viability theory [2].

2 PRELIMINARIES

For a fairly detailed introduction to reinforcement learning, the reader is referred to Appendix 8. In what follows, we discuss the architecture of Deep Q-Network.

2.1 Deep Q-Network (DQN)

Since a DQN is essentially an artificial neural network, or simply a neural network (NN), we begin by describing one. In particular, we discuss the architecture of a *fully connected feedforward network* with real-valued vector inputs. *Activation functions* form the basic building blocks of an NN. The typical domain for an activation function σ is \mathbb{R} , and its range \mathcal{R} is usually a subset of \mathbb{R} , i.e., $\sigma : \mathbb{R} \rightarrow \mathcal{R} \subset \mathbb{R}$. Depending on whether the range of σ , \mathcal{R} , is compact or unbounded, it is said to be *squashing* or *non-squashing*, respectively. There are many activation functions, the following are a few examples considered in this paper: (a) Sigmoid $[1/(1+e^{-x})]$, (b) Hyperbolic Tangent $[e^x - e^{-x}/(e^x + e^{-x})]$, (c) Gaussian Error Linear Unit $\left[x \int_{-\infty}^x e^{-y^2/2}/\sqrt{2\pi} dy \right]$, and (d) Sigmoid Linear Unit $[x/(1+e^{-x})]$.

An NN is a collection of *activations* that are arranged in a sequence of *layers*, starting with an *input* layer, then followed by one or more *hidden* layers, and ending with the *output* layer. An NN with two or more hidden layers is called a *Deep Neural Network* (DNN). Figure 2 illustrates one such NN architecture. By convention, an NN is constructed from left to right starting with the input layer and ending with the output layer. Further, the layers are arranged in a feedforward architecture, in that any two successive layers constitute a *complete bipartite graph* with edges directed from the left layer into the right.

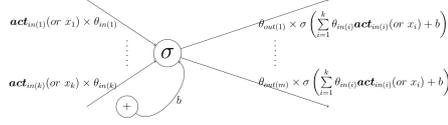


Figure 1: Single activation from some layer

Figure 1 illustrates a single activation σ within some layer. There are k edges leading into and m leading out of σ , where $m, k \geq 1$. When σ is in the input layer, the in-edges connect the k components of the input vector to it. As a part of other layers, the *in-edges* connect the k activation-outputs from the previous layer to its input. Further, each in-edge is associated with a weight that equals the product of the corresponding previous layer activation output $\mathbf{act}_{in(i)}$ (or input component x_i) and network-weight $\theta_{in(i)}$, $1 \leq i \leq k$. The input value to the activation is given by

$$\sum_{i=1}^k \mathbf{act}_{in(i)} \theta_{in(i)} + b$$

or $\sum_{i=1}^k x_i \theta_{in(i)} + b$, where b is a tunable bias term. Suppose σ is part of an input or hidden layer, then the edges leading out of it, the *out-edges*, connect its output

$$\sigma \left(\sum_{i=1}^k \mathbf{act}_{in(i)} \theta_{in(i)} + b \right) \quad (1)$$

(or $\sigma(\sum_{i=1}^k x_i \theta_{in(i)} + b)$) to the input of the m activations in the following layer. Finally, if σ is part of the output layer, its output (1) is combined with the output from other activations that also belong to the outer layer, to obtain the required NN output. For more details the reader may refer to [11,22].

Note on tunable biases: Subsequently, we assume that there are no tunable biases added to the activation inputs. In particular, we assume that the input is merely $\sum \theta_{in(i)} \mathbf{act}_{in(i)}$ (or $\sum \theta_{in(i)} x_i$ if the activation belongs to the input layer). We make this simplification for the sake of clarity in presentation. Our analysis will remain unaltered, except for minor bookkeeping, if one wishes to account for tunable biases.

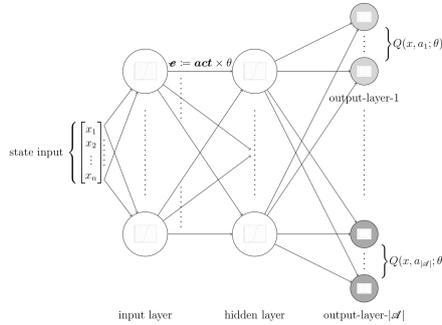


Figure 2: Schematic Representation of a DQN

We are now ready to discuss the DQN architecture, also illustrated in Fig. 2. Its input is the state vector $x \in \mathbb{S}$, and its output is a vector of dimension $|\mathcal{A}|$. The DQN output layer is a union of $|\mathcal{A}|$ separate (sub) output layers, one for each action. The output-layer- i associated with action a_i , is fully connected to the previous hidden layer, see Fig. 2. In particular, they are connected to the same layer. Let $l(a)$ be the number of activations in the output layer associated with action a , then $Q(x, a; \theta) := \sum_{i=1}^{l(a)} \mathbf{act}_{a(i)} \theta_{a(i)}$, where $\mathbf{act}_{a(i)}$ is the activation- i output and $\theta_{a(i)}$ is the associated network weight. *Note that we use $\theta_{a(i)}$ and $\mathbf{act}_{a(i)}$, instead of merely using θ_i and \mathbf{act}_i , respectively, to emphasize the association with action.*

In a nutshell, DQN is a parameterization of the vector $(Q^*(x, a))_{a \in \mathcal{A}}$, where Q^* is the optimal Q-function. In Deep Q-Learning, one updates the DQN weights $\theta := (\theta_e \mid e \text{ is an edge in the DQN})$ iteratively, in order to find θ^* such that $Q(x, a; \theta^*) \approx Q^*(x, a), \forall (x, a) \in \mathbb{S} \times \mathcal{A}$.

3 DEEP Q-LEARNING

To minimize the squared Bellman loss, Deep Q-learning iterates the update

$$\theta_{n+1} \leftarrow \theta_n + \gamma(n) \nabla_{\theta} \ell(\theta_n, x_n, a_n) \quad (2)$$

of the DQN weight vector $\theta \in \mathbb{R}^d$, where the following notation is used:

- (i) $\theta_n \in \mathbb{R}^d$, $x_n \in \mathbb{S}$, and $a_n \in \mathcal{A}$ for $n \in \mathbb{N}_0$. The state space \mathbb{S} is assumed to be \mathbb{R}^n for some $n \geq 1$, and \mathcal{A} is a finite set of actions.
- (ii) The loss gradient of (2) is given by

$$\begin{aligned} \nabla_{\theta} \ell(\theta_n, x_n, a_n) &= \nabla_{\theta} Q(x_n, a_n; \theta_n) \times \\ &\quad (r(x_n, a_n) + \alpha \max_{a' \in \mathcal{A}} Q(x_{n+1}, a'; \theta_n) - Q(x_n, a_n; \theta_n)), \end{aligned} \quad (3)$$

where α is the discount factor. Since a_n is the action taken at time n , $\nabla_{\theta} \ell(\theta_n, x_n, a_n)$ denotes the loss-gradient back-propagated via a_n .

- (iii) $\gamma(n)$ is the step size sequence satisfying the standard assumptions of non-summability and square summability.

Note that the loss gradient is calculated using the *sample value* $\max_{a' \in \mathcal{A}} Q(x_{n+1}, a'; \theta_n)$ instead of the *expected value* $\int \max_{a' \in \mathcal{A}} Q(x', a'; \theta_n) p(dx' \mid x, a)$. This is because the transition kernel p is unknown in real applications. The algorithm *observes* the next state x_{n+1} and the reward $r(x_n, a_n)$, after applying a_n in state x_n .

The state Markov process is determined by the transition kernel $p(dy \mid x, a)$. In training, actions are picked through a policy that *exploits* the approximation capability of DQN, while simultaneously *exploring* new actions. In other words, the transition kernel is indirectly influenced by the network weights. Hence, we denote the controlled transition kernel by $p(dy \mid x, a, \theta)$. For fixed weights θ and a fixed stochastic policy π_{θ} , the transition kernel is given by

$$\tilde{p}_{\theta}(dy \mid x) = \sum_{a \in \mathcal{A}} p(dy \mid x, a, \theta) \pi_{\theta}(x, da).$$

The policy is subscripted with θ to emphasize that it depends on the network weights (via exploitation). Let us suppose that π_θ only exploits and does not explore. Then the above stochastic policy is the Dirac measure given by $\pi_\theta(x, da) = \delta_{\underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(x, a; \theta)}$. Furthermore, the above kernel becomes

$$\tilde{p}_\theta(dy | x) = p\left(dy \mid x, \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(x, a; \theta), \theta\right).$$

Note on notation: We use dy and da (instead of just y and a) to represent the variables on which p and π_θ , respectively, define distributions, so as to easily distinguish the variable under consideration. Suppose π_θ is a Dirac measure. Then, through a slight abuse of notation, we use $\pi_\theta(x)$ to represent $\underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(x, a; \theta)$.

3.1 Assumptions

The assumptions required to analyze (2) are as follows:

- (A1) $\gamma(n) > 0$ for all $n \geq 0$, $\sum_{n \geq 0} \gamma(n) = \infty$ and $\sum_{n \geq 0} \gamma(n)^2 < \infty$. Further, the sequence monotonically decreasing.
- (A2) (a) $\sup_{n \geq 0} \|\theta_n\|_2 < \infty$ a.s., (b) $\sup_{n \geq 0} \|x_n\|_2 < \infty$ a.s.
- (A3) The state transition kernel $p(\cdot | x, a)$ is continuous in the x -coordinate.
- (A4) The DQN is composed of activation functions that are squashing and twice continuously differentiable.
- (A5) The reward function $r : \mathbb{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is continuous.

The first assumption regarding the step size sequence (learning rate) is standard in the literature. Recall that the loss gradient in (2) is calculated using samples that are supposed to approximate expected values. The resulting sampling errors are controlled using step sizes that are square summable. The stability assumption (A2) is essential for analyzing the long-term behavior of (2).

Consider two different but “closely neighbored” states in the environment. Assumptions (A3) and (A5) state that the consequences (successor states and rewards, respectively) of taking the same action in these states are similar. These assumptions are not only natural, but also ensure the performance of approximation-based algorithms like Deep Q-learning. As long as the state-action pairs encountered during training are a rich enough representation of $\mathbb{S} \times \mathcal{A}$, (A3) and (A5) facilitate good approximation of the Q-function.

The assumption of squashing activations (A4) is mainly made for the sake of clarity of presentation and can easily be relaxed. An extension to general (twice continuously differentiable) activations is provided in Section 5.

3.2 Properties of the loss gradient

The aim of this section is to prove certain useful properties that facilitate an abstract view of the loss gradient, with lesser “moving parts”. In particular, we show that $\nabla_\theta \ell$ is (A) locally Lipschitz continuous in the θ -coordinate, and (B)

continuous in the x and a -coordinates. Suppose we equip \mathcal{A} with the discrete topology. Then, since \mathcal{A} is a finite set, the resulting discrete space is compact, so that $\nabla_{\theta}\ell$ is trivially continuous in the a -coordinate. As for the rest, we relegate a couple of technical lemmata to Appendix 9, and summarize the required results in Lemma 1 below.

Let us define the sequence $\{M_n\}_{n \geq 0}$ as follows:

$$M_n := \sum_{m=0}^{n-1} \gamma(m)\psi_m, \quad n \geq 0, \quad \text{where}$$

$$\psi_m := \alpha \left[\max_{a \in \mathcal{A}} Q(x_{m+1}, a; \theta_m) - \int \max_{a \in \mathcal{A}} Q(x, a; \theta_m) p(dx \mid x_m, a_m, \theta_m) \right] \nabla_{\theta} Q(x_m, a_m; \theta_m).$$

It can be shown that $\{M_n\}_{n \geq 0}$ is a zero-mean martingale with respect to the filtration $\mathcal{F}_{n-1} := \sigma\langle x_m, a_m, \theta_m \mid m \leq n \rangle$, $n \geq 1$. Recall that we assume stability of (2) and the state sequence, i.e., $\sup_{n \geq 0} \|\theta_n\| < \infty$ and $\sup_{n \geq 0} \|x_n\| < \infty$ a.s. This, together with the twice continuous differentiability of Q in the θ -coordinate (shown in Lemma 9, Appendix 9), lets us conclude that $\sup_{n \geq 0} |Q(x_n, a_n; \theta_n)| < K_1 < \infty$, and that $\|\nabla_{\theta} Q(x_n, a_n; \theta_n)\| < K_2 < \infty$, where K_1 and K_2 are possibly sample-path dependent. Hence $\sup_{n \geq 0} \|\psi_n\| \leq K < \infty$, where K may again be sample-path dependent. Finally, the square summability of the step size sequence, assumption (A1), implies that $\sum_{m=0}^n \gamma(m)^2 \|M_m\|^2 < \infty$ a.s. Convergence of the martingale sequence $\{M_n\}_{n \geq 0}$ follows from the martingale convergence theorem, see [10].

Recall the loss gradient $\nabla_{\theta}\ell(\theta_n, x_n, a_n)$ given by (3), and let us rewrite it using the definition of ψ_n as

$$\begin{aligned} \nabla_{\theta}\ell(\theta_n, x_n, a_n) &= (r(x_n, a_n) + \\ &\quad \alpha \int \max_{a' \in \mathcal{A}} Q(y, a'; \theta_n) p(dy \mid x_n, a_n, \theta_n) \\ &\quad - Q(x_n, a_n; \theta_n)) \nabla_{\theta} Q(x_n, a_n; \theta_n) + \psi_n. \end{aligned} \quad (4)$$

Hence, (2) becomes

$$\theta_{n+1} \leftarrow \theta_n + \gamma(n) \left[\nabla_{\theta} \hat{\ell}(\theta_n, x_n, a_n) + \psi_n \right], \quad (5)$$

$$\begin{aligned} \text{where } \nabla_{\theta} \hat{\ell}(\theta_n, x_n, a_n) &:= (r(x_n, a_n) + \\ &\quad \alpha \int \max_{a' \in \mathcal{A}} Q(y, a'; \theta_n) p(dy \mid x_n, a_n, \theta_n) \\ &\quad - Q(x_n, a_n; \theta_n)) \nabla_{\theta} Q(x_n, a_n; \theta_n). \end{aligned}$$

Since the martingale sequence $\{M_n\}_{n \geq 0}$ converges a.s., the impact of ψ_n vanishes asymptotically. In other words, (2) and (5) are asymptotically identical to (have the same limiting set as)

$$\theta_{n+1} \leftarrow \theta_n + \gamma(n) \left[\nabla_{\theta} \hat{\ell}(\theta_n, x_n, a_n) \right]. \quad (6)$$

Note on notation: Rather than keeping track of two versions of the loss gradients, $\nabla_{\theta}\ell$ and $\nabla_{\theta}\hat{\ell}$ from equations (2) and (6), respectively, we redefine $\nabla_{\theta}\ell :=$

$\nabla_{\theta} \hat{\ell}$. With this slight abuse of notation, we hope to avoid unnecessary confusion. The reader does not need to track two different losses. In our subsequent analysis, when we refer to (2), the associated loss gradient is

$$\begin{aligned} \nabla_{\theta} \ell(\theta_n, x_n, a_n) &:= (r(x_n, a_n) + \alpha \\ &\int \max_{a' \in \mathcal{A}} Q(y, a'; \theta_n) p(dy | x_n, a_n, \theta_n) - \\ &Q(x_n, a_n; \theta_n)) \nabla_{\theta} Q(x_n, a_n; \theta_n). \end{aligned} \quad (7)$$

Lemma 1. $\nabla_{\theta} \ell(\theta_n, x_n, a_n)$, redefined as (7), is continuous and locally Lipschitz continuous in the θ -coordinate.

Proof. For the proof, one can combine the consequences of (i) Lemmas 8, 9 and 10 (see Appendix 9), (ii) assumption (A5), i.e., the continuity of the reward function r , and (iii) the fact that the sum and product of continuous and locally Lipschitz continuous functions are also continuous and locally Lipschitz continuous, respectively. \square

The Lipschitz constant from the above statement is local and changes with θ . However, as discussed before, following the proof of Lemma 9 (see Appendix 9), it also depends on x . If the domain of a locally Lipschitz continuous function is restricted to a compact subset, then the restricted function is Lipschitz continuous. Assumption (A2) states that $\sup_{n \geq 0} \|\theta_n\|_2 < \infty$ and $\sup_{n \geq 0} \|x_n\|_2 < \infty$ a.s. This can be used to conclude that $\nabla_{\theta} \ell$ is Lipschitz continuous in the θ -coordinate, when restricted to an appropriate compact subset of $\mathbb{R}^d \times \mathbb{S}$. We note that the Lipschitz constant may be sample-path dependent and refer to the proof of Lemma 1 in [17], where something very similar is shown.

4 CONVERGENCE ANALYSIS

To analyze the long-term behavior of (2), we first construct an associated continuous-time trajectory with identical limiting behavior. Then, instead of (2), we may analyze the continuous-time trajectory.

First, we divide the time axis $[0, \infty)$ using the given step size sequence as follows:

$$t_n = 0 \text{ and } t_n = \sum_{m=0}^{n-1} \gamma(m) \text{ for } n \geq 1.$$

We now define the required trajectory $\bar{\theta} \in C([0, \infty), \mathbb{R}^d)$ as follows:

- (a) $\bar{\theta}(t_n) = \bar{\theta}_n$, $n \geq 0$,
- (b) $\bar{\theta}(t) = \bar{\theta}(t_n) + \frac{t-t_n}{t_{n+1}-t_n} [\bar{\theta}(t_{n+1}) - \bar{\theta}(t_n)]$ for $t \in (t_n, t_{n+1})$ and $n \geq 0$.

As the sequence of actions taken are directly linked to the DQN-weights θ via ‘‘exploitation’’, we need to better understand them. To this end, we define the following measure process:

$$\mu(t) = \delta_{(x_n, a_n)}, \quad t \in [t_n, t_{n+1}),$$

where $\delta_{(x,a)}$ is the Dirac measure that places mass 1 on the state-action pair $(x, a) \in \mathbb{S} \times \mathcal{A}$. Hence $\mu : [0, \infty) \rightarrow \mathcal{P}(\mathbb{S} \times \mathcal{A})$ defines a process of probability

measures on $\mathbb{S} \times \mathcal{A}$. For our analysis, we need to define limits for the “left-shifted” measure process $\{\mu([t_n, \infty))\}_{n \geq 0}$. For that purpose, we first define a metric space (similar to the one from [5]) consisting of such measure processes below.

To start with, we observe that the action space \mathcal{A} is *compact metrizable*, as it is discrete and finite. As for \mathbb{S} , recall our assumption $\mathbb{S} = \mathbb{R}^n$. Thus, it follows from the Alexandroff extension that \mathbb{S} is one-point compactifiable. In particular, the inverse stereographic projection $S^{-1} : \mathbb{S} \rightarrow \mathcal{S}^n$ is such that $\mathcal{S}^n \setminus S^{-1}(\mathbb{S}) = (0, \dots, 0, 1)$, where \mathcal{S}^n represents the $(n+1)$ -dimensional Hausdorff compact sphere of radius 1 centered at the origin, and $(0, \dots, 0, 1)$ is the “north pole”. In other words, the inverse stereographic projection is the required compactification embedding of \mathbb{S} into \mathcal{S}^n , see [16].

Every measure $\nu \in \mathcal{P}(\mathbb{S} \times \mathcal{A})$ has a push forward counterpart in $\mathcal{P}(\mathcal{S}^n \times \mathcal{A})$. It places mass 0 on $(0, \dots, 0, 1) \times \mathcal{A}$. *Moving forward, note that we shall use the same symbol to represent both the measure and its push forward counterpart.* Also note that $\mathcal{S}^n \times \mathcal{A}$ is *compact Hausdorff in the product topology*.

Let us define \mathcal{U} to be the space of all measurable functions $\nu(\cdot) = \nu(\cdot, dx, da)$ from $[0, \infty)$ to $\mathcal{P}(\mathcal{S}^n \times \mathcal{A})$.

Lemma 2. *\mathcal{U} is compact metrizable. Further, this metric coincides with the coarsest topology that renders continuous the map*

$$\nu \mapsto \int_0^T g(t) \int f d\nu(t) dt,$$

for all, $T > 0$, $f \in \mathbb{C}(\mathcal{S}^n \times \mathcal{A})$ and $g \in \mathbb{L}^2([0, T], \mathbb{R})$.

Proof. By emulating the proof of Lemma 3 in [5] with “ $\mathcal{S}^n \times \mathcal{A}$ ” replacing “ $\bar{\mathbb{S}}$ ”, and making appropriate modifications, the required proof is obtained. We do not repeat it here, to avoid redundancies. \square

Define $\tilde{\nabla} \ell(\theta, \nu) := \int \nabla_{\theta} \ell(\theta, x, a) \nu(dx, da)$, where $\nu \in \mathcal{P}(\mathbb{S}, \mathcal{A})$. Lemma 1 implies that $\tilde{\nabla} \ell$ is *continuous in both coordinates and locally Lipschitz continuous in the θ -coordinate*. Further, $\|\tilde{\nabla} \ell(\theta, \nu)\| \leq K(1 + \|\theta\|)$, i.e., its growth is bounded as a function of θ alone. Let us also define the following sequence of trajectories in $\mathbb{C}([0, \infty), \mathbb{R}^d)$: $\theta^n(t) = \bar{\theta}(t_n) + \int_0^t \tilde{\nabla} \ell(\theta^n(s), \mu^n(s)) ds$, where $\mu^n(t) := \mu(t_n + t)$, $t \geq 0$ and $n \geq 0$. In other words, we consider solutions to the set of non-autonomous ordinary differential equations: $\left\{ \dot{\theta}^n(t) = \tilde{\nabla} \ell(\theta^n(t), \mu^n(t)) \right\}_{n \geq 0}$. As stated earlier, to understand the long-term behavior of (2), one can study the behavior of the limit of sequence $\{\bar{\theta}([t_n, \infty))\}_{n \geq 0}$, in $\mathbb{C}([0, \infty), \mathbb{R}^d)$ as $n \rightarrow \infty$. We can show the following property.

Lemma 3. *For every $T > 0$,*

$$\lim_{n \rightarrow \infty} \sup_{t \in [0, T]} \|\bar{\theta}(t_n + t) - \theta^n(t)\| = 0 .$$

Proof. Please refer to Appendix 10.1 for a proof. \square

Then, instead of (2) or the associated trajectory $\bar{\theta}$, we could focus on the sequence of trajectories $\{\theta^n([0, \infty))\}_{n \geq 0}$. Now we may tap into the rich literature of tools and techniques available from viability theory [2, 3].

The family of trajectories $\{\theta^n([0, \infty))\}_{n \geq 0}$, in $\mathbb{C}([0, \infty), \mathbb{R}^d)$, is equicontinuous and point-wise bounded. It follows from the Arzela-Ascoli theorem [4] that it is sequentially compact. Note that the topology of $\mathbb{C}([0, \infty), \mathbb{R}^d)$ is the one induced by the topologies of $\mathbb{C}([0, T], \mathbb{R}^d)$ for every $0 < T < \infty$. Now, let us consider the family $\{\mu^n\}_{n \geq 0} \subset \mathcal{U}$. As \mathcal{U} is a compact metric space, $\{\mu^n\}_{n \geq 0}$ is sequentially compact. Hence, there is a common subsequence $\{m(n)\} \subset \{n\}$ such that $\mu^{m(n)} \rightarrow \mu^\infty$ in \mathcal{U} and $\theta^{m(n)} \rightarrow \theta^\infty$ in $\mathbb{C}([0, \infty), \mathbb{R}^d)$. *With a slight abuse of notation, we have $\mu^n \rightarrow \mu^\infty$ in \mathcal{U} and $\theta^n \rightarrow \theta^\infty$ in $\mathbb{C}([0, \infty), \mathbb{R}^d)$.* In other words, the sequences μ^n and θ^n are convergent in their respective spaces.

Below we state another important result, namely that convergence of the measure process in \mathcal{U} implies convergence in distribution of the corresponding measure sequence, at every point in time.

Lemma 4. *If $\mu^n \rightarrow \mu^\infty$ in \mathcal{U} , then a.e. $\mu^n(t) \rightarrow \mu^\infty(t)$ in $\mathcal{P}(\mathbb{S} \times \mathcal{A})$ for $t \in [0, \infty)$.*

Proof. We begin by recalling that the same notation is used to denote a measure on $\mathbb{S} \times \mathcal{A}$ and its push forward counterpart on $\mathcal{S}^n \times \mathcal{A}$. It follows from the definition of convergence in \mathcal{U} that $\int_0^T g(s) \int f \mu^n(s, dx, da) ds \rightarrow \int_0^T g(s) \int f \mu(s, dx, da) ds$ as $n \rightarrow \infty$, for every $g \in \mathbb{L}^2([0, T], \mathbb{R})$ and $f \in \mathbb{C}(\mathcal{S}^n \times \mathcal{A})$. We claim that this implies, for every $f \in \mathbb{C}(\mathcal{S}^n \times \mathcal{A})$, $\int f \mu^n(s, dx, da) \rightarrow \int f \mu(s, dx, da)$ a.e. for $s \in [0, \infty)$. Once this claim is proven, we can conclude that *s-a.e.* $\mu^n(s) \rightarrow \mu^\infty(s)$ in $\mathcal{P}(\mathcal{S}^n \times \mathcal{A})$, which finally yields the lemma.

To prove the claim, let us assume the contrary. In particular, we assume $\exists f \in \mathbb{C}(\mathcal{S}^n \times \mathcal{A})$, $T > 0$, $\epsilon > 0$ and a non-zero Lebesgue measure set $A \in \mathcal{B}([0, T])$, such that at least one of the following properties holds for all $s \in A$:

- (a) $\liminf_{n \rightarrow \infty} \int f \mu^n(s, dx, da) - \int f \mu(s, dx, da) > \epsilon$,
- (b) $\liminf_{n \rightarrow \infty} \int f \mu^n(s, dx, da) - \int f \mu(s, dx, da) < -\epsilon$,
- (c) $\limsup_{n \rightarrow \infty} \int f \mu^n(s, dx, da) - \int f \mu(s, dx, da) > \epsilon$,
- (d) $\limsup_{n \rightarrow \infty} \int f \mu^n(s, dx, da) - \int f \mu(s, dx, da) < -\epsilon$.

We only present arguments for case (a), as the corresponding ones for the others are identical. Since f is bounded, we apply the Dominated Convergence Theorem (DCT) [10] to conclude that

$$\begin{aligned} \liminf_{n \rightarrow \infty} \int_0^T \mathbb{1}_A \left[\int f \mu^n(s, dx, da) - \int f \mu(s, dx, da) \right] ds \\ > \epsilon l(A) > 0, \end{aligned}$$

where $l(A)$ denotes the Lebesgue measure of A . This directly contradicts the definition of convergence of measures in \mathcal{U} .

It is left to show that $\mu^n(t) \rightarrow \mu^\infty(t)$ in $\mathcal{P}(\mathbb{S} \times \mathcal{A})$ a.e. for $t \in [0, \infty)$. To do this, we pick $t \in [0, \infty)$ such that $\mu^n(t) \rightarrow \mu^\infty(t)$ in $\mathcal{P}(\mathcal{S}^n \times \mathcal{A})$ and show that their pull back versions converge in $\mathcal{P}(\mathbb{S} \times \mathcal{A})$. This is done by showing that

$\limsup_{n \rightarrow \infty} \mu^n(t, C) \leq \mu^\infty(t, C)$ for every closed set $C \in \mathcal{B}(\mathbb{S} \times \mathcal{A})$ (Portmanteau theorem [4]).

We first observe that the measures $\{\mu^n(t)\}_{0 \leq n \leq \infty}$ are tight as a consequence of (A2). Hence they place a mass of 0 on $(0, \dots, 0, 1) \times \mathcal{A}$. If we restrict these measures to $S^{-1}(\mathbb{S}) \times \mathcal{A}$, then $\mu^n|_{S^{-1}(\mathbb{S}) \times \mathcal{A}} \xrightarrow{d} \mu^\infty|_{S^{-1}(\mathbb{S}) \times \mathcal{A}}$. Next, we consider an arbitrary closed subset $C \in \mathcal{B}(\mathbb{S} \times \mathcal{A})$. Since the stereographic projection is bicontinuous, $\hat{C} := \{(S^{-1}(x), a) \mid (x, a) \in C\}$ is closed in $S^{-1}(\mathbb{S}) \times \mathcal{A}$, equipped with subspace topology (with respect to $\mathcal{S}^n \times \mathcal{A}$). Clearly, $\limsup_{n \rightarrow \infty} \mu^n(t, \hat{C}) \leq \mu^\infty(t, \hat{C})$. Now, as $\mu^n(t, \hat{C})$ is the push forward measure of $\mu^n(t, C)$ for all $0 \leq n \leq \infty$, we obtain the required result. \square

We can use one of the many available measurable selection theorems [19] to drop the a.e. clause in the statement of Lemma 4. Hence, we have hitherto shown that $\theta^n \rightarrow \theta^\infty$ in $\mathbb{C}([0, \infty), \mathbb{R}^d)$ and $\mu^n(s) \xrightarrow{d} \mu^\infty(s)$ for all $s \in [0, \infty)$. We now need to show that θ^∞ is a solution of $\dot{\theta}(t) = \tilde{\nabla} \ell(\theta(t), \mu^\infty(t))$. Then, one can study the limiting behavior of a solution to the ODE $\dot{\theta}(t) = \tilde{\nabla} \ell(\theta(t), \mu^\infty(t))$, to understand the long-term behavior of Deep Q-Learning given by (2).

Lemma 5. θ^∞ is a solution to $\dot{\theta}(t) = \tilde{\nabla} \ell(\theta(t), \mu^\infty(t))$.

Proof. Fix an arbitrary $T > 0$. We need to show that

$$\sup_{t \in [0, T]} \left\| \theta^n(t) - \theta^\infty(0) - \int_0^t \tilde{\nabla} \ell(\theta^\infty(s), \mu^\infty(s)) ds \right\| \rightarrow 0.$$

Let us first consider the following:

$$\left\| \theta^n(0) + \int_0^t \tilde{\nabla} \ell(\theta^n(s), \mu^n(s)) ds - \theta^\infty(0) - \int_0^t \tilde{\nabla} \ell(\theta^\infty(s), \mu^\infty(s)) ds \right\|, \quad (8)$$

$$\begin{aligned} & \|\theta^n(0) - \theta^\infty(0)\| + \left\| \int_0^t \tilde{\nabla} \ell(\theta^n(s), \mu^n(s)) ds - \int_0^t \tilde{\nabla} \ell(\theta^\infty(s), \mu^n(s)) ds - \int_0^t \tilde{\nabla} \ell(\theta^\infty(s), \mu^\infty(s)) ds \right\|. \quad (9) \end{aligned}$$

Next, we note the following:

- (A) From Lemma 4 we have $\mu^n(s) \xrightarrow{d} \mu^\infty(s)$ (converges in distribution on $\mathbb{S} \times \mathcal{A}$) for all $s \in [0, T]$.

- (B) From (A2), i.e., the stability of the algorithm, and the boundedness of $\nabla_\theta \ell$ as a function of θ , we get $\nabla_\theta \ell(\theta^\infty(s), \cdot) \in \mathbb{C}_b(\mathbb{S} \times \mathcal{A})$. Hence, as a consequence of note (A), $\int \nabla_\theta \ell(\theta^\infty(s), x, a) \mu^n(s) \rightarrow \int \nabla_\theta \ell(\theta^\infty(s), x, a) \mu^\infty(s)$ for all $s \in [0, T]$.

Using DCT, we get

$$\left\| \int_0^t \tilde{\nabla} \ell(\theta^\infty(s), \mu^n(s)) ds - \int_0^t \tilde{\nabla} \ell(\theta^\infty(s), \mu^\infty(s)) ds \right\| \rightarrow 0. \quad (10)$$

Further, it follows from the Arzela-Ascoli theorem that the convergence in (10) is uniform over $[0, T]$.

Since $\tilde{\nabla} \ell$ is locally Lipschitz continuous in θ , we get

$$\begin{aligned} & \left\| \int_0^t \tilde{\nabla} \ell(\theta^n(s), \mu^n(s)) ds - \int_0^t \tilde{\nabla} \ell(\theta^\infty(s), \mu^n(s)) ds \right\| \\ & \leq L \int_0^t \|\theta^n(s) - \theta^\infty(s)\| ds. \end{aligned} \quad (11)$$

As $\theta^n \rightarrow \theta^\infty$ uniformly over $[0, T]$, the l.h.s. of (11) $\rightarrow 0$ uniformly over $[0, T]$. The discussion surrounding (10) and (11) implies that (9) $\rightarrow 0$ and hence (8) $\rightarrow 0$, uniformly over $[0, T]$. As T is *arbitrary*, the lemma follows. \square

To develop a better understanding of Deep Q-Learning, we need to study μ^∞ , the limiting distribution over the state-action pairs. In the following lemma, we show that $\mu^\infty(t, dx \times \mathcal{A})$ is stationary with respect to the state Markov process, $\forall t \geq 0$. Recall that $p(\cdot | x, a, \theta)$ is the controlled transition kernel of the state Markov process. We use $p(\cdot | x, \mathcal{A}, \theta)$ to denote the probability associated with transitioning out of state x (when some action is picked). We use $p(dy | x, \mathcal{A}, \theta) \mu(dx \times \mathcal{A})$ to denote $\int_{\mathcal{A}} p(dy | x, a, \theta) \mu(dx, da)$. In words, it represents the probability to transition from state x to state y , given that $(x, a) \sim \mu$.

Lemma 6. *For all $t \in [0, \infty)$, $\mu^\infty(t, dy \times \mathcal{A}) = \int_{\mathbb{S}} p(dy | x, \mathcal{A}, \theta^\infty(t)) \mu^\infty(t, dx \times \mathcal{A})$. In other words, the limiting marginal constitutes a stationary distribution over the state Markov process. Further, $\{\mu^\infty(t, dx, da)\}_{t \geq 0}$ is tight.*

For a proof of this lemma, we refer to Appendix 10.2.

Tightness of $\{\mu^\infty(t, dx, da)\}_{t \geq 0}$ implies that it is relative compact in the Prokhorov metric. This property, combined with the stability of (2), yields $\{n(k)\}_{k \geq 0} \subset \{n\}_{n \geq 0}$, such that both $\lim_{n(k) \rightarrow \infty} \bar{\theta}(t_{n(k)})$ and $\lim_{n(k) \rightarrow \infty} \mu(t_{n(k)}, dx, da)$ have limits in \mathbb{R}^d and $\mathcal{P}(\mathbb{S} \times \mathcal{A})$, respectively. The properties of these limits, let us call them $\bar{\theta}^\infty$ and $\bar{\mu}^\infty$, determine the long-term behavior of (2). Lemmas 8 to 6 were stated and proved to build up to the most important result of this paper, which concerns the limiting behavior of (2). We state and prove this result below, followed by a discussion of its implications.

Theorem 1. *Assuming (A1)–(A5), the limit $\bar{\theta}^\infty$ of the deep Q-learning algorithm, i.e., iteration (2), is such that $\tilde{\nabla}\ell(\bar{\theta}^\infty, \bar{\mu}^\infty) = 0$ and $\bar{\mu}^\infty(dx \times \mathcal{A})$ is a stationary distribution of the state Markov process x .*

Proof. From previous lemmas we know that (2) tracks θ , a solution to the non-autonomous ODE $\dot{\theta}(t) = \tilde{\nabla}\ell(\theta(t), \mu^\infty(t))$. Further, there is a sample path dependent compact subset of \mathbb{R}^d , \mathcal{K} , such that θ remains inside of it. This is because the algorithm is assumed to be stable, i.e., $\theta_n \in \mathcal{K} \forall n \geq 0$. To determine the limit of the algorithm, $\bar{\theta}^\infty$, we need $\lim_{t \rightarrow \infty} \theta(t)$.

To analyze $\dot{\theta}(t) = \tilde{\nabla}\ell(\theta(t), \mu^\infty(t))$, we transform it into an autonomous ODE through the standard change of variables trick. For this, we define $s(t) := \frac{t}{1+t}$, then $\dot{s}(t) = (1-s(t))^2$ and $t = \frac{s(t)}{1-s(t)}$. We get the following transformed autonomous ODE:

$$\begin{aligned} & (\dot{\theta}(t), \dot{s}(t)) = \\ & \left(\tilde{\nabla}\ell\left(\theta(t), \mu^\infty\left(\frac{s(t)}{1-s(t)}\right)\right), (1-s(t))^2 \right). \end{aligned} \quad (12)$$

Before proceeding, we state the following useful theorem, paraphrased to suit our purpose:

[Theorem 2, Chapter 6 of [2]] *Let F be a continuous map from a closed subset $\hat{\mathcal{X}} \subset \mathcal{X}$ to \mathcal{X} . Let $x(\cdot)$ be a solution trajectory of $\dot{x}(t) = F(x(t))$, such that it is inside $\hat{\mathcal{X}}$. Then, the solution converges to x^* , an equilibrium of F .*

To utilize the theorem, we define the following: $\mathcal{X} := \mathbb{R}^d \times [0, 1]$, $\hat{\mathcal{X}} := \mathcal{K} \times [0, 1]$, and $F: \hat{\mathcal{X}} \rightarrow \mathcal{X}$ such that $F(\theta, s) := \left(\tilde{\nabla}\ell\left(\theta, \mu^\infty\left(\frac{s}{1-s}\right)\right), (1-s)^2 \right)$. It now follows from the above theorem that the transformed ODE (12) converges to $(\bar{\theta}^\infty, 1)$, an equilibrium of F . Further, 1 is the unique equilibrium point of $(1-s)^2$, and $\bar{\theta}^\infty$ is an equilibrium of $\tilde{\nabla}\ell(\bar{\theta}^\infty, \bar{\mu}^\infty)$, where $\lim_{t \rightarrow \infty} \mu^\infty(t) \xrightarrow{d} \bar{\mu}^\infty$. We discussed the existence of the limit $\bar{\mu}^\infty$ in the paragraph before stating this theorem.

Lemma 6 shows that $\mu^\infty(t)$ is a stationary distribution of the state Markov process x for all $t \geq 0$, i.e.,

$$\mu^\infty(t, dy \times \mathcal{A}) = \int_{\mathcal{S}} p(dy | x, \mathcal{A}, \theta^\infty(t)) \mu^\infty(t, dx \times \mathcal{A}).$$

Letting $t \rightarrow \infty$ on both sides of the above equation yields

$$\bar{\mu}^\infty(dy \times \mathcal{A}) = \int_{\mathcal{S}} p(dy | x, \mathcal{A}, \bar{\theta}^\infty) \bar{\mu}^\infty(dx \times \mathcal{A}).$$

In other words, the marginal over the states, $\bar{\mu}^\infty(dx \times \mathcal{A})$, is stationary with respect to the state process. \square

4.1 On practical implications of the theory

The primary goal of Deep Q-Learning is to find the optimal DQN-weights θ^* such that $\operatorname{argmax}_{a \in \mathcal{A}} Q(x, a; \theta^*) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(x, a)$, where Q^* is the optimal Q-function. This is achieved by minimizing the squared Bellman loss.

Theorem 1 states that the Deep Q-Learning algorithm given by (2) converges to $\bar{\theta}^\infty$, a local minimizer of the average squared Bellman loss. The averaging over state-action pairs is induced by the limiting measure $\bar{\mu}^\infty \in \mathcal{P}(\mathbb{S} \times \mathcal{A})$. In particular, we have

$$\int \nabla_{\theta} \ell(\bar{\theta}^\infty, x, a) \bar{\mu}^\infty(dx, da) = 0. \quad (13)$$

Lemma 6 states that the limiting marginal distribution $\bar{\mu}^\infty(dx \times \mathcal{A})$, over the state space \mathbb{S} , is stationary. Deep Q-Learning is typically employed in complex environments with multiple stationary distributions. Since $\bar{\mu}^\infty$ captures the long-term behavior of the training process, it directly depends on the distribution of the data encountered during training. As the squared Bellman loss is minimized on average in accordance to $\bar{\mu}^\infty$, the quality of learning is entirely captured by $\bar{\mu}^\infty$. In particular, the trained DQN approximates the optimal Q-factors accurately for state-action pairs that are distributed in accordance to $\bar{\mu}^\infty$. Performance is therefore good when encountering states arising from the “limiting marginal”.

Fix $a \in \mathcal{A}$ and let $\mathbb{S}(a)$ be a measurable subset of \mathbb{S} such that a is the optimal action associated with every $x \in \mathbb{S}(a)$. For the sake of illustration, we consider a scenario wherein $\bar{\mu}^\infty(\mathbb{S}(a) \times \mathcal{A}) > 0$ and $\bar{\mu}^\infty(\mathbb{S}(a) \times a) = 0$. Roughly speaking, the set of state-action pairs given by $\{(x, a) \mid x \in \mathbb{S}\}$ were not encountered during training. This could happen, for example, due to poor exploration-exploitation trade-offs, or due to improper initialization of the DQN weights. The Q-factors may hence be poorly approximated on $\mathbb{S}(a) \times a$, and the trained DQN-agent cannot be expected to take optimal actions in these states. This explains the observation that, in practice, Deep Q-Learning sometimes fails to generalize well beyond the data encountered during training. Existing literature (see e.g. [21, 23]) does not account for such behaviors. Since DQN is usually trained using a simulator, it may be possible to empirically estimate $\bar{\mu}^\infty$. This knowledge may help identify scenarios wherein DQN is undertrained, thereby avoiding circumstances like the one sketched above.

5 Weakening (A4) to allow twice continuously differentiable non-squashing activation functions

The hitherto presented analysis accounts for DQN architectures with differentiable squashing activations. In this section, we discuss modifications to our analysis that allow for general activations as well. In particular, the modifications account for activations such as Sigmoid Linear Unit (SiLU), Gaussian Error Linear Unit (GELU), etc.

Let us begin by understanding the role of squashing activations in our analysis. In Lemma 8, the squashing property is used to find a x -independent \hat{L} such that $|Q(x, a; \theta)| \leq \hat{L} \|\theta\|_2$. Note that Lemma 8 is true even when the activations are non-squashing, provided \mathbb{S} is a compact metric space. Since (A2) states that $\sup_{n \geq 0} \|x_n\|_2 < \infty$ a.s., there is a sample path dependent compact set $\mathbb{S}_c \subset \mathbb{S}$ such that $x_n \in \mathbb{S}_c \forall n \geq 0$. Using this information, we may modify the statement of Lemma 8 as follows:

Lemma 7. $\forall \theta \in \mathbb{R}^d \sup_{a \in \mathcal{A}} |Q(x, a; \theta)| \leq \tilde{L} \|\theta\|_2$, and $\tilde{L} > 0$ is dependent on x . Further, there is a sample path dependent \hat{L} , independent of x , such that $\sup_{x \in \mathbb{S}_c} \sup_{a \in \mathcal{A}} |Q(x, a; \theta)| \leq \hat{L} \|\theta\|_2$, where \mathbb{S}_c is as defined above.

Parts of the analysis using Lemma 8 must now be modified to use Lemma 7. Other Lemmata, for e.g., Lemma 10 do not change when using Lemma 7 instead of Lemma 8.

6 Extension to account for experience replay

Now, we extend our analysis to account for experience replay, an idea that allows the RL agent to relearn from past experiences. Specifically, at time T , the agent has ready access to $\{(x_k, a_k, r(x_k, a_k), x_{k+1})\}_{T-H+1 \leq k \leq T}$, the history of states encountered, actions taken, rewards received and transitions made. The optimal size of the experience replay H is problem dependent, and tunable. At time T , to update the NN weights θ , the agent first samples a mini-batch of size $\hat{H} < H$ from the experience replay and calculates the following average loss gradient:

$$\frac{1}{\hat{H}} \sum_{i=1}^{\hat{H}} \nabla_{\theta} \ell(\theta_T, x_{k(T,i)}, a_{k(T,i)}), \text{ where}$$

$$T - H + 1 \leq k(T, i) \leq T.$$

The DQN weights are updated as follows:

$$\theta_{n+1} = \theta_n + \gamma(n) \left[\frac{1}{\hat{H}} \sum_{i=1}^{\hat{H}} \nabla_{\theta} \ell(\theta_n, x_{k(n,i)}, a_{k(n,i)}) \right]. \quad (14)$$

To analyze (14), we must redefine μ . For $t \in [t_n, t_{n+1})$, redefine $\mu(t)$ to be the probability measure (on $\mathbb{S} \times \mathcal{A}$) that places a mass of $1/\hat{H}$ on $(x_{k(n,i)}, a_{k(n,i)})$ for $1 \leq i \leq \hat{H}$. With the new definition of μ , for $t = t_n$ we get:

$$\tilde{\nabla} \ell(\bar{\theta}(t), \mu(t)) = \int \nabla_{\theta} \ell(\bar{\theta}(t), x, a) \mu(t) =$$

$$\frac{1}{\hat{H}} \sum_{i=1}^{\hat{H}} \nabla_{\theta} \ell(\theta_n, x_{k(n,i)}, a_{k(n,i)}).$$

Emulating the proofs of the Lemmata up to Lemma 5 for the new μ , shows that (14) tracks a solution to the non-autonomous o.d.e. $\dot{\theta}(t) = \tilde{\nabla} \ell(\theta(t), \mu^{\infty}(t))$. Again, μ^{∞} is a limit of the redefined measure process sequence $\{\mu([t, \infty))\}_{t \geq 0}$ in \mathcal{U} .

Lemma 6 states the the limiting marginal measure process $\mu^{\infty}(t, dx \times \mathcal{A})$ is stationary with respect to the state Markov process for every $t \geq 0$. For it to hold in the presence of experience replay we redefine ξ_n and \mathcal{F}_n as follows:

$$\xi_n := \sum_{m=0}^{n-1} \frac{1}{\hat{H}} \left[\sum_{i=1}^{\hat{H}} (f(x_{k(m,i)+1}) - \int f(y) p(dy | x_{k(m,i)}, a_{k(m,i)}, \theta_{k(m,i)})) \right],$$

$\mathcal{F}_{n-1} = \sigma(x_m, a_m, \theta_m, \Xi_m \mid m \leq n)$ for $n \geq 1$, where $\{\Xi_n\}_{n \geq 0}$ is the random process associated with mini-batch sampling. Typically the mini-batches are all sampled independently over time, hence $\{\Xi_n\}_{n \geq 0}$ constitutes an independent sequence of random variables. With these modifications the rest the steps involved in the proof of Lemma 6 may be readily emulated. This would directly lead to the statement of the main result, Theorem 1. In conclusion, Deep Q-Learning with experience replay, (14), converges to $\hat{\theta}^\infty$ such that $\nabla_{\theta} \ell(\hat{\theta}^\infty, \hat{\mu}^\infty) = 0$, where $\hat{\mu}^\infty$ is a limit of $\{\tilde{\mu}^\infty(t)\}_{t \geq 0}$ as $t \rightarrow \infty$, and $\tilde{\mu}^\infty$ is the limiting measure process of the redefined μ -process. Again, $\hat{\mu}^\infty(dx \times \mathcal{A})$ is stationary with respect to the state Markov process.

It is a common belief among deep learning practitioners that experience replay plays an important role in stabilizing the DQN training. In regards to the long-term behavior, we show that the use of experience replay has a qualitative effect on learning. *This is because the limiting measure $\tilde{\mu}^\infty$ is shaped by the mini-batches sampled from experience replay during training, and it is richer than the one resulting from no experience replay.*

7 CONCLUSION

In this paper, we presented an asymptotic analysis of Deep Q-Learning under practical and verifiable assumptions. An important contribution is the complete characterization of the DQN performance as a function of training. We obtained this result by analyzing the limit of a closely associated measure process (on the state-action pairs). The result has various implications that we shall elaborate on more closely in future work. In particular, it helps explain empirical observations regarding the performance of Deep Q-Learning that current theory does not account for. Practically motivated extensions and generalizations like this one are also on our agenda of future work.

References

- [1] Joshua Achiam, Ethan Knight, and Pieter Abbeel. Towards characterizing divergence in deep q-learning. *arXiv preprint arXiv:1903.08894*, 2019.
- [2] J-P Aubin and Arrigo Cellina. *Differential inclusions: set-valued maps and viability theory*, volume 264. Springer Science & Business Media, 2012.
- [3] Jean-Pierre Aubin, Alexandre M Bayen, and Patrick Saint-Pierre. *Viability theory: new directions*. Springer Science & Business Media, 2011.
- [4] Patrick Billingsley. *Convergence of probability measures*. John Wiley & Sons, 2013.
- [5] Vivek S. Borkar. Stochastic approximation with ‘controlled Markov’ noise. *Systems & control letters*, 55(2):139–145, 2006.
- [6] Vivek S Borkar. *Stochastic approximation: a dynamical systems viewpoint*, volume 48. Springer, 2009.
- [7] Qi Cai, Zhuoran Yang, Chi Jin, and Zhaoran Wang. Provably efficient exploration in policy optimization. *arXiv preprint arXiv:1912.05830*, 2019.

- [8] Zaiwei Chen, Sheng Zhang, Thinh T Doan, John-Paul Clarke, and Siva Theja Maguluri. Finite-sample analysis of nonlinear stochastic approximation with applications in reinforcement learning. *arXiv e-prints*, pages arXiv-1905, 2019.
- [9] Simon S Du, Jason D Lee, Gaurav Mahajan, and Ruosong Wang. Agnostic q -learning with function approximation in deterministic systems: Near-optimal bounds on approximation error and sample complexity. *Advances in Neural Information Processing Systems*, 33, 2020.
- [10] Rick Durrett. *Probability: Theory and Examples*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2010.
- [11] Daniel Graupe. *Principles of artificial neural networks*, volume 7. World Scientific, 2013.
- [12] Seungchan Kim, Kavosh Asadi, Michael Littman, and George Konidaris. Deepmellow: removing the need for a target network in deep q-learning. In *Proceedings of the Twenty Eighth International Joint Conference on Artificial Intelligence*, 2019.
- [13] Harold Kushner and G George Yin. *Stochastic approximation and recursive algorithms and applications*, volume 35. Springer Science & Business Media, 2003.
- [14] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *ICLR (Poster)*, 2016.
- [15] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [16] James R. Munkres. *Topology. 2nd ed.* Upper Saddle River, NJ: Prentice Hall, 2nd ed. edition, 2000.
- [17] Arunselvan Ramaswamy, Adrian Redder, and Daniel E Quevedo. Optimization over time-varying networks with unbounded delays. *arXiv preprint arXiv:1912.07055*, 2019.
- [18] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
- [19] Daniel H Wagner. Survey of measurable selection theorems. *SIAM Journal on Control and Optimization*, 15(5):859–903, 1977.
- [20] Pan Xu and Quanquan Gu. A finite-time analysis of q-learning with neural network function approximation. In *International Conference on Machine Learning*, pages 10555–10565. PMLR, 2020.
- [21] Zhuora Yang, Yuchen Xie, and Zhaoran Wang. A theoretical analysis of deep q-learning. *arXiv preprint arXiv:1901.00137*, 2019.

- [22] Bayya Yegnanarayana. *Artificial neural networks*. PHI Learning Pvt. Ltd., 2009.
- [23] Shaofeng Zou, Tengyu Xu, and Yingbin Liang. Finite-sample analysis for sarsa with linear function approximation. In *Advances in Neural Information Processing Systems*, pages 8668–8678, 2019.

8 Appendix: Preliminaries: Reinforcement Learning (RL)

In RL an *agent* interacts with an *environment* over time, via *actions*. It takes the current (environment) *state* into consideration to pick an action, and receives a feedback in terms of a *reward*. The environment then moves to a new state. This is schematically represented in Figure 3. *The goal in RL is to ensure that the agent takes a sequence of actions, such that the rewards accumulated over time are maximized.*

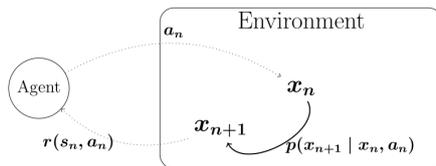


Figure 3: Snapshot of interaction at step n

Formally speaking, the above stated interactions can be modelled as a Markov Decision Process (MDP). It is defined as a 5-tuple $(\mathbb{S}, \mathcal{A}, p, r, \alpha)$, where:

\mathbb{S} is the state space. In typical applications $\mathbb{S} \equiv \mathbb{R}^k$, $k > 0$.

\mathcal{A} is the action space. In this paper, \mathcal{A} is a discrete finite set.

p is the “controlled” transition kernel. We use $p(\cdot | x, a)$ to represent the distribution of the next state given the current state and action.

r is the reward function. In particular, $r(x, a)$ denotes the reward associated with taking action a at state x .

α is the discount factor with $0 < \alpha \leq 1$. It is used to discount the relevance of future consequences of actions.

A policy π is defined as a function from \mathbb{S} to \mathcal{A} . Given π , we can associate a *Value function* $V^\pi(x)$ with each $x \in \mathbb{S}$, with $V^\pi(x) := \mathbb{E} \left[\sum_{n \geq 0} \alpha^n r(x_n, \pi(x_n)) \mid x_0 = x \right]$.

The goal in RL can be restated to find π^* such that $V^{\pi^*}(x) = \max_{\pi} V^\pi(x)$ for all $x \in \mathbb{S}$. In Dynamic Programming parlance π^* is a solution to the *infinite horizon discounted reward problem*.

Closely related to the value function is the concept of Q-function, defined over state-action pairs $(x, a) \in \mathbb{S} \times \mathcal{A}$ by $Q^\pi(x, a) := r(x, a) + \alpha \int V^\pi(x') p(dx' | x, a)$, where π is a fixed policy. The optimal Q-function is defined as:

$$Q^*(x, a) := r(x, a) + \alpha \int V^{\pi^*}(x') p(dx' | x, a).$$

Clearly, $\max_{a \in \mathcal{A}} Q^*(x, a) = V^{\pi^*}(x)$ and $\pi^*(x) = \operatorname{argmax}_{a \in \mathcal{A}} Q^*(x, a)$ for all $x \in \mathbb{S}$. Hence, in order to find π^* it is sufficient to find Q^* . This is the idea behind Q-Learning. Its variant, Deep Q-Learning, has shown tremendous promise in solving complex problems involving continuous state spaces, where Q-Learning typically fails. It involves parameterizing the optimal Q-function using a DNN, called the Deep Q-Network (DQN). The goal is to find the optimal set of parameters (DQN weights) θ^* , *by interacting with the environment*, such that $Q(x, a; \theta^*) \approx Q^*(x, a)$ for all $(x, a) \in \mathbb{S} \times \mathcal{A}$. The DQN is trained to minimize the following squared Bellman loss over all state-action pairs (x, a) :

$$\left[r(x, a) + \alpha \int \max_{a' \in \mathcal{A}} Q(x', a'; \theta) p(dx' | x, a) - Q(x, a; \theta) \right]^2.$$

9 Appendix: Technical lemmas supporting Lemma 1

Let us recall that every action is associated with a different output layer:

$$Q(x, a; \theta) = \sum_{i=1}^{l(a)} \mathbf{act}_a(i) \theta_a(i), \text{ with } l(a) \text{ the width of the layer associated with action } a.$$

Lemma 8. $\sup_{x \in \mathbb{S}, a \in \mathcal{A}} |Q(x, a; \theta)| \leq \hat{L} \|\theta\|_2$, for some $\hat{L} > 0$.

Proof. We begin by noting that activation functions considered herein are also squashing. Hence, absolute values of their outputs are bounded by some $0 < c < \infty$. Let us fix arbitrary $x \in \mathbb{S}$ and $a \in \mathcal{A}$, then

$$|Q(x, a; \theta)| \leq c \sum_{i=1}^{l(a)} |\theta_a(i)| = c \|\theta_a\|_1,$$

where $\|\cdot\|_1$ is the 1-norm. It now follows from $\|\theta_a\|_1 \leq l(a) \|\theta_a\|_2$, that $|Q(x, a; \theta)| \leq cl(a) \|\theta_a\|_2$. If we let $\hat{L} := c l(a)$, then the statement of the lemma follows. \square

Since we allow for possibly unbounded Q-factors, the above lemma indicates that we need arbitrarily large DQN weights for good approximation. Depending on the system states encountered during training, the Deep Q-Learning algorithm explores an appropriate subspace associated with the weight vector. Hence, the approximation capability of the trained DQN depends on the state-action pairs encountered during training. The difference, in state distributions, between the training and test scenarios will determine performance.

Recall that we parameterize the Q-function using a neural network that consists of twice continuously differentiable activation functions. Hence, Q may be viewed as a composition of twice continuously differentiable activations, and the DQN weight vector. In other words, Q itself is twice continuously differentiable. This intuition is formalized in the next lemma.

Lemma 9. $Q(x, a; \theta)$ is twice continuously differentiable in the θ -coordinate for every $x \in \mathbb{S}$ and $a \in \mathcal{A}$, where θ is the DQN weight vector.

Proof. Recall that the DQN weights are updated using the back propagation algorithm, i.e., the chain rule. Given the DQN weight-vector $\theta \in \mathbb{R}^d$, we need to show that $\partial^2 Q(\hat{x}, \hat{a}; \theta) / \partial \theta_i^2$ exists and is continuous for $1 \leq i \leq d$. Also, recall from the *note on tunable biases* at the end of Section 2.1, that without loss of generality we may only consider tunable edge weights, and ignore tunable bias terms.

Let us fix an arbitrary $\hat{x} \in \mathbb{S}$, $\hat{a} \in \mathcal{A}$ and $i \in \{1, \dots, d\}$. DQN weight θ_i is associated with an edge of the NN. Also associated with this edge is another weight $\mathbf{e}_i := \mathbf{act}_i \theta_i$, where \mathbf{act}_i is the output of an activation from the previous layer (from the head of the edge). This is illustrated in Fig. 4.

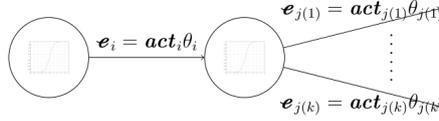


Figure 4: Section of a DNN

To prove the lemma, we show something stronger, i.e., that both $\partial^2 Q(\hat{x}, \hat{a}; \theta) / \partial \theta_i^2$ and $\partial^2 Q(\hat{x}, \hat{a}; \theta) / \partial \mathbf{e}^2$ are continuous. The proof involves inducting on the depth of the DNN, starting from the output layer, and going backwards. Note that $Q(\hat{x}, \hat{a}; \theta) = \sum_{i=1}^{l(\hat{a})} \mathbf{e}_{\hat{a}(i)}$, where $\mathbf{e}_{\hat{a}(i)} := \mathbf{act}_{\hat{a}(i)} \theta_{\hat{a}(i)}$, where $l(\hat{a})$ is the number of activations in the output layer of action \hat{a} . Also, note that $\mathbf{act}_{\hat{a}(i)}$ is the output of the i -th activation in the output layer associated with action \hat{a} , and $\theta_{\hat{a}(i)}$ is the corresponding network edge-weight, $1 \leq i \leq l(\hat{a})$, see Section 2.1 for details. We have, $\partial^2 Q(\hat{x}, \hat{a}; \theta) / \partial \theta_{a(i)}^2 = \partial^2 Q(\hat{x}, \hat{a}; \theta) / \partial \mathbf{e}_{a(i)}^2 = \mathbf{0}$ for all $a \neq \hat{a}$, where *subscript $a(i)$ is used to indicate that $\theta_{a(i)}$ and $\mathbf{e}_{a(i)}$ are associated with the output layer of action a* . Twice continuous differentiability with respect to $\theta_{\hat{a}(i)}$ and $\mathbf{e}_{\hat{a}(i)}$ directly follows from the same property of the activation units, $1 \leq i \leq l(\hat{a})$.

Let us assume that the hypothesis is true for weights associated with edges out of the $(l+1)^{st}$ layer and prove for the l^{th} layer. Fig. 4 illustrates an edge out of an l^{th} layer activation, and its associated weight $\mathbf{e}_i := \mathbf{act}_i \theta_i$, where $i \in \{1, \dots, d\}$. Also note that, in the Fig. 4, $\mathbf{act}_{j(m)} = \mathbf{act}_0$ for all $1 \leq m \leq k$. It follows directly from the back-propagation algorithm (chain rule) that:

$$\begin{aligned} \frac{\partial Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_i} &= \frac{\partial \mathbf{act}_0}{\partial \mathbf{e}_i} \sum_{m=1}^k \left[\frac{\partial Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_{j(m)}} \theta_{j(m)} \right], \\ \frac{\partial^2 Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_i^2} &= \left(\frac{\partial \mathbf{act}_0}{\partial \mathbf{e}_i} \right)^2 \sum_{m=1}^k \left[\frac{\partial^2 Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_{j(m)}^2} \theta_{j(m)}^2 \right] + \\ &\quad \frac{\partial^2 \mathbf{act}_0}{\partial \mathbf{e}_i^2} \sum_{m=1}^k \left[\frac{\partial Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_{j(m)}} \theta_{j(m)} \right]. \end{aligned}$$

From the induction hypothesis and the twice continuous differentiability of $\hat{\mathbf{act}}_0$, we get that $\frac{\partial^2 Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_i^2}$ is continuous. Next, we observe:

$$\begin{aligned} \frac{\partial Q(\hat{x}, \hat{a}; \theta)}{\partial \theta_i} &= \frac{\partial Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_i} \frac{\partial \mathbf{e}_i}{\partial \theta_i} = \mathbf{act}_i \frac{\partial Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_i}, \\ \frac{\partial^2 Q(\hat{x}, \hat{a}; \theta)}{\partial \theta_i^2} &= (\mathbf{act}_i)^2 \frac{\partial^2 Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_i^2}. \end{aligned}$$

The continuity of $\frac{\partial^2 Q(\hat{x}, \hat{a}; \theta)}{\partial \theta_i^2}$ follows from the twice continuous differentiability of $\frac{\partial^2 Q(\hat{x}, \hat{a}; \theta)}{\partial \mathbf{e}_i^2}$. \square

Since Q is two times continuously differentiable in the θ -coordinate, it is locally Lipschitz continuous in that coordinate. Also, the Lipschitz constant may depend on x , in addition to θ . Let us fix arbitrary $\hat{a} \in \mathcal{A}$ and $\hat{\theta} \in \mathbb{R}^d$. Since $Q(\cdot, \hat{a}; \hat{\theta})$ and $\nabla_{\theta} Q(\cdot, \hat{a}; \hat{\theta})$ are composed (via addition and multiplication) of twice continuously differentiable functions (activation units), we get that both Q and $\nabla_{\theta} Q$ are continuous in the x -coordinate. Although we do not need it here, the stronger property of local Lipschitz continuity may also be shown. Finally, note that Q and $\nabla_{\theta} Q$ are continuous in the a -coordinate, since \mathcal{A} is finite.

Lemma 10. *The following map is continuous and locally Lipschitz continuous in the θ -coordinate:*

$$(x, a, \theta) \mapsto \int \max_{a \in \mathcal{A}} Q(x', a; \theta) p(dx' | x, a, \theta).$$

Proof. We begin by fixing arbitrary $\hat{x} \in \mathbb{S}$ and $\hat{a} \in \mathcal{A}$. Given $\theta \in \mathbb{R}^d$, Lemma 9 implies the existence of $\mathcal{N}(\theta, \hat{x})$, without loss of generality a compact neighborhood of θ , and $L(\theta, \hat{x}) > 0$, such that $\forall \theta_1, \theta_2 \in \mathcal{N}(\theta, \hat{x})$:

$$|Q(\hat{x}, \hat{a}; \theta_1) - Q(\hat{x}, \hat{a}; \theta_2)| \leq L(\theta, \hat{x}) \|\theta_1 - \theta_2\|_2.$$

Since \hat{a} is fixed, $p(\cdot | \hat{x}, \hat{a}, \theta) \equiv p(\cdot | \hat{x}, \hat{a})$, i.e., the transition kernel does not depend on θ . Recall that the dependence of p on θ is only via the action a . Define $a_1(x) := \operatorname{argmax}_{a \in \mathcal{A}} Q(x, a; \theta_1)$, then following the above line of thought (with “ x ” replacing “ \hat{x} ” and “ $a_1(x)$ ” replacing “ \hat{a} ”) we get:

$$\begin{aligned} & \left| \max_{a \in \mathcal{A}} Q(x, a; \theta_1) - \max_{a \in \mathcal{A}} Q(x, a; \theta_2) \right| \\ & \leq |Q(x, a_1(x); \theta_1) - Q(x, a_1(x); \theta_2)| \\ & \leq L(\theta, x) \|\theta_1 - \theta_2\|. \end{aligned} \tag{15}$$

Hence, from Lemma 8 and the compactness of $\mathcal{N}(\theta, \hat{x})$, we conclude that

$$\sup_{\hat{\theta} \in \mathcal{N}(\theta, \hat{x})} \sup_{x \in \mathbb{S}} \sup_{a \in \mathcal{A}} |Q(x, a; \hat{\theta})| < \infty.$$

In particular, there exists a bounded measurable function $\hat{F}_\theta : x \mapsto L(x, \theta)$ such that (15) is satisfied for every $x \in \mathbb{S}$, with $\hat{F}_\theta(x)$ as the Lipschitz constant. Hitherto presented arguments and observations yield:

$$\begin{aligned} & \left| \int \max_{a \in \mathcal{A}} Q(x, a; \theta_1) p(dx \mid \hat{x}, \hat{a}, \theta_1) - \right. \\ & \quad \left. \int \max_{a \in \mathcal{A}} Q(x, a; \theta_2) p(dx \mid \hat{x}, \hat{a}, \theta_2) \right| \leq \\ & \leq \|\theta_1 - \theta_2\|_2 \int L(\theta, x) p(dx \mid \hat{x}, \hat{a}) \leq L \|\theta_1 - \theta_2\|_2, \end{aligned} \quad (16)$$

where $L = 2 \times \sup_{\hat{\theta} \in \mathcal{N}(\theta, \hat{x})} \sup_{x \in \mathbb{S}} \sup_{a \in \mathcal{A}} |Q(x, a; \hat{\theta})|$.

Let us fix arbitrary $\hat{\theta} \in \mathbb{R}^d$ and $\hat{a} \in \mathcal{A}$. Define $\hat{a}(x) \in \operatorname{argmax}_{a \in \mathcal{A}} Q(x, a; \hat{\theta})$

and $\hat{Q}(x) := Q(x, \hat{a}(x), \hat{\theta})$ for all $x \in \mathbb{S}$. Note that there may be many actions that maximize the Q function, $\hat{a}(\cdot)$ selecting one of them. First, we show that $x_n \rightarrow x$ implies that $\hat{Q}(x_n) \rightarrow \hat{Q}(x)$, and hence that $\hat{Q} \in \mathbb{C}_b(\mathbb{S})$ (from Lemma 8). To this end, we show that every subsequence of $\{\hat{Q}(x_n)\}_{n \geq 0}$ has a further subsequence that converges, and the limit always equals $\hat{Q}(x)$. Let us begin by considering the entire sequence itself. Since \mathcal{A} is a compact metric space, $\exists \{n(m)\}_{m \geq 0} \subset \{n\}_{n \geq 0}$ such that $\hat{a}(x_{n(m)}) \rightarrow \hat{a}$ for some $\hat{a} \in \mathcal{A}$, hence $Q(x_{n(m)}, \hat{a}(x_{n(m)}); \hat{\theta}) \rightarrow Q(x, \hat{a}; \hat{\theta})$. We claim that $\hat{a} = \hat{a}(x)$, thus implying $\hat{Q}(x_{n(m)}) \rightarrow \hat{Q}(x)$. To see that the claim is true assume the contrary. In other words, $\hat{a}(x) \neq \hat{a}$ and $Q(x, \hat{a}(x); \hat{\theta}) > Q(x, \hat{a}; \hat{\theta}) + \epsilon$, for some $\epsilon > 0$. From the continuity of Q , we get that $\exists M > 0$ with $|Q(x_{n(m)}, \hat{a}(x); \hat{\theta}) - Q(x, \hat{a}(x); \hat{\theta})| \leq \epsilon/4$ and $|Q(x_{n(m)}, \hat{a}(x_{n(m)}); \hat{\theta}) - Q(x, \hat{a}; \hat{\theta})| \leq \epsilon/4$, for all $m \geq M$. Hence, we get that $Q(x_{n(m)}, \hat{a}(x); \hat{\theta}) > Q(x_{n(m)}, \hat{a}(x_{n(m)}); \hat{\theta})$, a contradiction. Finally, we note that the above set of arguments can be repeated starting with any subsequence of $\{n\}_{n \geq 0}$.

Now that we have $\hat{Q} \in \mathbb{C}_b(\mathbb{S})$, we are ready to prove continuity in the x -coordinate. Recall that we have assumed the transition kernel to be continuous in x . Hence $x_n \rightarrow x$ implies that $p(\cdot \mid x_n, \hat{a}, \hat{\theta}) \xrightarrow{d} p(\cdot \mid x, \hat{a}, \hat{\theta})$, i.e., the kernels converge in distribution. It now follows from the definition of ‘‘convergence in distribution’’ that $\int \hat{Q}(y) p(dy \mid x_n, \hat{a}, \hat{\theta}) \rightarrow \int \hat{Q}(y) p(dy \mid x, \hat{a}, \hat{\theta})$. In other words, we have the required, namely

$$\begin{aligned} x_n \rightarrow x & \implies \int \max_{a \in \mathcal{A}} Q(y, a, \hat{\theta}) p(dy \mid x_n, \hat{a}, \hat{\theta}) \rightarrow \\ & \int \max_{a \in \mathcal{A}} Q(y, a, \hat{\theta}) p(dy \mid x, \hat{a}, \hat{\theta}) \end{aligned}$$

as $n \rightarrow \infty$. Finally, recall that \mathcal{A} is compact metrizable as it is a finite. Hence continuity in the a -coordinate is trivial. \square

10 Appendix: Missing Proofs

10.1 Proof of Lemma 3

Proof. First we define the notation $[t]$ for $t \geq 0$ as $[t] := t_{\sup\{n|t_n \leq t\}}$. Next, we need to show that:

$$\sup_{t \in [0, T]} \|\bar{\theta}(t_n + t) - \bar{\theta}([t_n + t])\| \in \Theta(\gamma(n)).$$

For this, we fix $t \in [0, T]$, then $[t_n + t] = t_{n+k}$ for some $k \geq 0$. Recall that

$$\bar{\theta}(t_n + t) = \bar{\theta}(t_{n+k}) + \frac{t_n + t - t_{n+k}}{\gamma(n+k)} (\bar{\theta}(t_{n+k+1}) - \bar{\theta}(t_{n+k})).$$

We use the following: $\|\bar{\theta}(t_{n+k+1}) - \bar{\theta}(t_{n+k})\| \leq \gamma(n+k) \|\nabla_{\theta} \ell(\bar{\theta}(t_{n+k}), x_{n+k}; a_{n+k})\|$; the stability of the algorithm, i.e., (A2); the monotonic property of the step-size sequence, i.e., (A1); and the boundedness of $\nabla_{\theta} \ell$ as a function of θ , to obtain $\|\bar{\theta}(t_n + t) - \bar{\theta}(t_{n+k})\| \in \Theta(\gamma(n))$. Similarly, let us show that:

$$\sup_{t \in [0, T]} \|\theta^n(t) - \theta^n([t_n + t] - t_n)\| \in \Theta(\gamma(n)).$$

Again, $[t_n + t] = t_{n+k}$ for some $k \geq 0$. We also have $\|\theta^n(t) - \theta^n(t_{n+k} - t_n)\| = \left\| \int_{t_{n+k} - t_n}^t \tilde{\nabla} \ell(\theta(s), \mu^n(s)) ds \right\|$. Using arguments similar to the ones made before, the required statement directly follows. It follows from all of the above arguments that it is enough to show the following in order to prove the lemma:

$$\sup_{t \in [0, T]} \|\bar{\theta}([t_n + t]) - \theta^n([t_n + t] - t_n)\| \rightarrow 0.$$

Once again we let $[t_n + t] = t_{n+k}$ for some $k \geq 0$, and observe that

$$\begin{aligned} & \|\bar{\theta}([t_n + t]) - \theta^n([t_n + t] - t_n)\| \leq \\ & \sum_{m=n}^{n+k-1} \int_{t_m}^{t_{m+1}} \|\tilde{\nabla} \ell(\bar{\theta}([s]), \mu^n(s - t_n)) - \\ & \quad \tilde{\nabla} \ell(\theta^n(s - t_n), \mu^n(s - t_n))\| ds, \\ & \|\bar{\theta}([t_n + t]) - \theta^n([t_n + t] - t_n)\| \leq \\ & \sum_{m=n}^{n+k-1} \int_{t_m}^{t_{m+1}} L \|\bar{\theta}([s]) - \theta^n(s - t_n)\|. \end{aligned}$$

Adding and subtracting $\theta^n([s] - t_n)$, the R.H.S. of above equation is less than or equal to

$$\begin{aligned} & \sum_{m=n}^{n+k-1} L \int_{t_m}^{t_{m+1}} \|\theta^n(s - t_n) - \theta^n([s] - t_n)\| + \\ & \sum_{m=n}^{n+k-1} L \int_{t_m}^{t_{m+1}} \|\bar{\theta}([s]) - \theta^n([s] - t_n)\|. \end{aligned}$$

Considering that $\|\bar{\theta}(t_n + t) - \bar{\theta}(t_{n+k})\|$ and $\|\theta^n(t) - \theta^n([t_n + t] - t_n)\| \in \Theta(\gamma(n))$, we get $\sum_{m=n}^{n+k-1} \int_{t_m}^{t_{m+1}} \|\theta^n(s - t_n) - \theta^n([s] - t_n)\| \leq \sum_{m=n}^{n+k-1} \Theta(\gamma(m)^2)$, which goes to zero as $n \rightarrow \infty$. Now we use the discrete version of Gronwall's inequality to get:

$$\begin{aligned} \|\bar{\theta}([t_n + t]) - \theta^n([t_n + t] - t_n)\| \leq \\ \left(L \sum_{m=n}^{n+k-1} \Theta(\gamma(m))^2 \right) \exp(LT). \end{aligned}$$

□

10.2 Proof of Lemma 6

Proof. Pick f from $\mathbb{C}_b(\mathbb{S})$, the convergence determining class for $\mathcal{P}(\mathbb{S})$. Without loss of generality, we assume that $0 \leq f \leq 1$. We define the following zero mean Martingale with respect to the filtration $\mathcal{F}_{n-1} := \sigma\langle x_m, a_m, \theta_m \mid m \leq n \rangle$, for $n \geq 1$:

$$\xi_n := \sum_{m=0}^{n-1} \gamma(m) \left[f(x_{m+1}) - \int_{\mathbb{S}} f(y) p(dy \mid x_m, a_m, \theta_m) \right]. \quad (17)$$

Since f is bounded and $\sum_{n \geq 0} \gamma(n)^2 < \infty$, the quadratic variation process associated with the above Martingale is convergent. It follows from the Martingale Convergence Theorem [10] that ξ_n converges almost surely. Hence for $t > 0$,

$$\sum_{m=n}^{\tau(n,t)} \gamma(m) \left[f(x_{m+1}) - \int_{\mathbb{S}} f(y) p(dy \mid x_m, a_m, \theta_m) \right] \rightarrow 0 \text{ a.s.}, \quad (18)$$

where $\tau(n,t) := \min\{m \geq n \mid t_m \geq t_n + t\}$. Since the steps-sizes are eventually decreasing, hence $\sum_{m=n}^{\tau(n,t)} [\gamma(m) - \gamma(m+1)] f(x_{m+1}) \rightarrow 0$ a.s. Then (18) becomes:

$$\sum_{m=n}^{\tau(n,t)} \gamma(m) \left[f(x_m) - \int_{\mathbb{S}} f(y) p(dy \mid x_m, a_m, \theta_m) \right] \rightarrow 0 \text{ a.s.} \quad (19)$$

Using the definition of μ , we rewrite (19) as:

$$\begin{aligned} \int_{t_n}^{t_n+t} \int_{\mathbb{S} \times \mathcal{A}} \left[f(x) - \int_{\mathbb{S}} f(y) p(dy \mid x, a, \bar{\theta}(s)) \right] \mu(s, dx, da) ds \\ \rightarrow 0 \text{ a.s.} \end{aligned} \quad (20)$$

Let us define a new function $\hat{f}(x, a) := f(x)$ for all $(x, a) \in \mathbb{S} \times \mathcal{A}$, then $\hat{f} \in \mathbb{C}_b(\mathbb{S} \times \mathcal{A})$. Since $\mu(t_n + \cdot) \rightarrow \mu^\infty(\cdot)$ in \mathcal{U} , it follows that as $n \rightarrow \infty$:

$$\begin{aligned} \int_{t_n}^{t_n+t} \int_{\mathbb{S} \times \mathcal{A}} \hat{f}(x, a) \mu(s, dx, da) ds \rightarrow \\ \int_0^t \int_{\mathbb{S} \times \mathcal{A}} \hat{f}(x, a) \mu^\infty(s, dx, da) ds. \end{aligned} \quad (21)$$

Further, the limit in (21) equals $\int_0^t \int_{\mathbb{S}} f(x) \mu^\infty(s, dx \times \mathcal{A}) ds$.

Recall that $(x, a, \theta) \mapsto p(\cdot | x, a, \theta)$ is a continuous map. Since f is a convergence determining function in $\mathcal{P}(\mathbb{S})$, it follows that $\int_{\mathbb{S}} f(y) p(dy | x, a, \bar{\theta}(s)) \rightarrow \int_{\mathbb{S}} f(y) p(dy | x, a, \theta^\infty(s))$ for all $s \in [0, t]$. Define $h_n(s, x, a) := \int_{\mathbb{S}} f(y) p(dy | x, a, \bar{\theta}(t_n + s))$ and $h_\infty(s, x, a) := \int_{\mathbb{S}} f(y) p(dy | x, a, \theta^\infty(s))$. For a fixed $s \in [0, t]$, $h_n(s, \cdot)$, $n \geq 0$, and $h_\infty(s, \cdot)$ belong to $\mathbb{C}_b(\mathbb{S} \times \mathcal{A})$. Hence,

$$\begin{aligned} \int_{\mathbb{S} \times \mathcal{A}} h_n(s, x, a) \mu(t_n + s, dx, da) &\rightarrow \\ \int_{\mathbb{S} \times \mathcal{A}} h_\infty(s, x, a) \mu^\infty(s, dx, da). \end{aligned} \quad (22)$$

It then follows from Dominated Convergence Theorem (DCT) [10] that:

$$\begin{aligned} \int_{t_n}^{t_n+t} \int_{\mathbb{S} \times \mathcal{A}} h_n(s, x, a) \mu(s, dx, da) ds &\rightarrow \\ \int_0^t \int_{\mathbb{S} \times \mathcal{A}} h_\infty(s, x, a) \mu^\infty(s, dx, da) ds. \end{aligned} \quad (23)$$

In other words, we have

$$\begin{aligned} \int_{t_n}^{t_n+t} \int_{\mathbb{S} \times \mathcal{A}} \int_{\mathbb{S}} f(y) p(dy | x, a, \bar{\theta}(s)) \mu(s, dx, da) ds &\rightarrow \\ \int_0^t \int_{\mathbb{S} \times \mathcal{A}} \int_{\mathbb{S}} f(y) p(dy | x, a, \theta^\infty(s)) \mu^\infty(s, dx, da) ds. \end{aligned} \quad (24)$$

From (20), (21) and (24) we get:

$$\begin{aligned} \int_0^t \int_{\mathbb{S} \times \mathcal{A}} f(x) \mu^\infty(s, dx, da) ds &= \\ \int_0^t \int_{\mathbb{S} \times \mathcal{A}} \int_{\mathbb{S}} f(y) p(dy | x, a, \theta^\infty(s)) \mu^\infty(s, dx, da) ds. \end{aligned} \quad (25)$$

Using Lebesgue's theorem we get that a.e. on $[0, t]$:

$$\begin{aligned} \int_{\mathbb{S} \times \mathcal{A}} f(x) \mu^\infty(s, dx, da) &= \\ \int_{\mathbb{S} \times \mathcal{A}} \int_{\mathbb{S}} f(y) p(dy | x, a, \theta^\infty(s)) \mu^\infty(s, dx, da). \end{aligned}$$

Applying Fubini's theorem [10] to swap the double integral on the R.H.S. of the above equation, gives us:

$$\begin{aligned} \int_{\mathbb{S}} f(x) \mu^\infty(s, dx, \mathcal{A}) &= \\ \int_{\mathbb{S}} f(y) \int_{\mathbb{S}} p(dy | x, \mathcal{A}, \theta^\infty(s)) \mu^\infty(s, dx, \mathcal{A}). \end{aligned}$$

Since f is a convergence determining function, we get that $\mu^\infty(s, dy, \mathcal{A}) = \int_{\mathbb{S}} p(dy | x, \mathcal{A}, \theta^\infty(s)) \mu^\infty(s, dx, \mathcal{A})$. Hence, we have shown that the limiting

distribution over the state-action pairs μ^∞ is such that, almost everywhere on $[0, \infty)$, its marginal over the state space constitutes a stationary distribution over the state Markov process with transition kernel $p(\cdot | x, \mathcal{A}, \theta)$.

Now, it is left to show that the family of measures $\{\mu^\infty(t, dx, da)\}_{t \geq 0}$ is tight. From previous discussions and observations, given $t \geq 0$, we can find $\{n(m)\}_{m \geq 0} \subset \{n\}_{n \geq 0}$ such that

$$\lim_{n(m) \rightarrow \infty} \mu(t_{n(m)}, dx, da) \xrightarrow{d} \mu^\infty(t, dx, da).$$

Using the Portmanteau Theorem [4], we get $\mu^\infty(t, \mathcal{K} \times \mathcal{A}') \geq \limsup_{n(m) \rightarrow \infty} \mu(t_{n(m)}, \mathcal{K} \times \mathcal{A}')$, where $\mathcal{K} \subset \mathbb{S}$ is compact and $\mathcal{A}' \subset \mathcal{A}$. Given $\epsilon > 0$, there exists $\mathcal{K}(\epsilon) \subset \mathbb{S}$, compact, such that $\inf_{m \geq 0} \mu(t_{n(m)}, \mathcal{K}(\epsilon) \times \mathcal{A}') \geq 1 - \epsilon$ for any $\mathcal{A}' \subset \mathcal{A}$, as $\mu(t_{n(m)})_{m \geq 0}$ is tight. Hence $\mu^\infty(t, \mathcal{K}(\epsilon) \times \mathcal{A}') \geq 1 - \epsilon$. As t was arbitrary, we get that $\{\mu^\infty(t, dx, da)\}_{t \geq 0}$ is tight. \square