How to Discount Deep Reinforcement Learning: Towards New Dynamic Strategies

Vincent François-Lavet
University of Liege
v.francois@ulq.ac.be

Raphael Fonteneau

University of Liege raphael.fonteneau@ulg.ac.be

Damien Ernst University of Liege dernst@ulg.ac.be

Abstract

Using deep neural nets as function approximator for reinforcement learning tasks have recently been shown to be very powerful for solving problems approaching real-world complexity such as [1]. Using these results as a benchmark, we discuss the role that the discount factor may play in the quality of the learning process of a deep Q-network (DQN). When the discount factor progressively increases up to its final value, we empirically show that it is possible to significantly reduce the number of learning steps. When used in conjunction with a varying learning rate, we empirically show that it outperforms original DQN on several experiments. We relate this phenomenon with the instabilities of neural networks when they are used in an approximate Dynamic Programming setting. We also describe the possibility to fall within a local optimum during the learning process, thus connecting our discussion with the exploration/exploitation dilemma.

1 Introduction

Reinforcement learning is a learning paradigm aiming at learning optimal behaviors while interacting within an environment [2]. One of the main challenge met when designing reinforcement learning algorithms is the fact that the state space may be very large or continuous, potentially leading to the fact that state(-action) value functions may not be represented comprehensively (for instance, using lookup tables) [3]. Recently, it has been empirically demonstrated in [1] that using deep neural networks as a function approximator may be very powerful in situations with high-dimensional sensory inputs such as the Atari games benchmark [4]. However, one of the main drawback of using (deep) neural networks as function approximator is that they may potentially become unstable when combined with a Q-Learning-type recursion [5].

In this paper, we propose to discuss the role that the discount factor may play in the stability and convergence of deep reinforcement learning algorithms. We empirically show that an increasing discount factor has the potential to improve the quality of the learning process. We also discuss the link between the discount factor influence, the instabilities of (deep) neural networks and the vanilla exploration/exploitation dilemma.

A motivation for this work comes from empirical studies about the attentional and cognitive mechanisms in delay of gratification. One well known experiment in the domain was a series of studies in which a child was offered a choice between one small reward provided immediately or two small rewards if they waited for a short period ("marshmallow experiment" [6]). The capacity to wait longer for the preferred rewards seems to develop markedly only at about ages 3-4. By 5 years old, most children are able to demonstrate better self-control by gaining control over their immediate

desires. According to this theory, it seems plausible that following immediate desires at a young age is a better way to develop its abilities and that delaying strategies is only advantageous afterwards when pursuing longer term goals. Similarly, reinforcement learning may also have an advantage of starting to learn by maximizing rewards on a short-term horizon and progressively giving more weights to delayed rewards.

The remaining of the paper is organized the following: we first recall the main equations used in the deep reinforcement learning problem formulation originally introduced in [1]. We then discuss the factors that influence the stability of (deep) neural networks. In light of this analysis, we suggest the possibility to modify the discount factor along the way to convergence up to its final value in order to speed up the learning process. To illustrate our approach, we use the benchmark proposed in [1]. In this context, we then discuss the role of the learning rate as well as the level of exploration.

2 Problem Formulation

We consider tasks in which an agent interacts with an environment so as to maximize expected future rewards

$$Q^*(s,a) = \max_{\pi} \mathbb{E}[r_t, \gamma r_{t+1}, \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi]$$
 (1)

which is the maximum sum of rewards r discounted by γ at each time-step t, achievable by a behaviour stochastic policy $\pi:S\times A\to [0,1]$ where $\pi(s,a)$ denotes the probability that action a may be chosen by policy π in state s. In this paper, the task is described by a time-invariant stochastic discrete-time system whose dynamics can be described by the following equation:

$$s_{t+1} = f(s_t, a_t, w_t)$$
 (2)

where for all t, s_t is an element of the state space S, the action a_t is an element of the action space A and the random disturbance w_t is an element of the disturbance space W generated by a time-invariant conditional probability distribution $w_t \sim P(.|s,a)$. In our experiments, the dynamics will be given by the Atari emulator, the state space will be based on observed pixels of the screen and the action space is the set of legal game actions $A = \{1, ..., K\}$.

The (unique) solution of the Bellman equation for the Q-value function [7] is given by:

$$Q^*(s,a) = (HQ^*)(s,a)$$
(3)

where H is an operator mapping any function $K: S \times A \to \mathbb{R}$ and defined as follows:

$$(HK)(s,a) = \mathbb{E}[r(s,a,w) + \gamma \max_{a' \in A} K(f(s,a',w),a')] \tag{4}$$

For most problems approaching real-world complexity, we need to learn a parameterized value function $Q(s,a;\theta_k)$. General function-approximation system such as neural networks are well suited to deal with high-dimensional sensory inputs. In addition, they work readily online, i.e. they can make use of additional samples obtained as learning happens by using only one supervised learning problem. In the case where a neural network $Q(s,a;\theta_k)$ is used to aim at convergence towards $Q^*(s,a)$, the parameters θ_k may be updated by stochastic gradient descent (or a variant), updating the current value $Q(s,a;\theta_k)$ towards a target value $Y_k^Q = r(s,a,w) + \gamma \arg\max_{a' \in A} Q(f(s,a',w),a';\theta_k^-)$ where θ_k^- refers to parameters from some previous Q-network. The Q-learning update when using the squared-loss amounts in updating the weights :

$$\theta_{k+1} = \theta_k + \alpha (Y_k^Q - Q(s, a; \theta_k)) \nabla_{\theta_k} Q(s, a; \theta_k)$$
(5)

where α is a scalar step size called the learning rate. A sketch of the algorithm is given in Figure 1.

3 Instabilities of the online neural fitted Q-learning

The Q-learning rule from Equation 5 can be directly implemented online using a neural network (as is the case for the Deep Q-Network). However, due to the generalization and extrapolation abilities of neural nets, they can build unpredictable changes at different places in the state-action space. It is known that errors may be propagated and that this may even become unstable. Additional care must therefore be taken since it can not be guaranteed that the current estimation for the accumulated

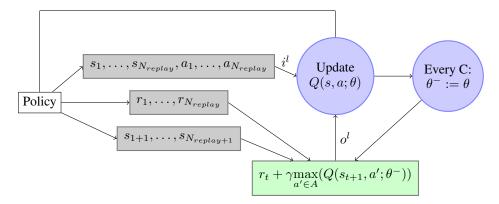


Figure 1: Sketch of the algorithm. $Q(s,a;\theta)$ is initialized to close to 0 everywhere on its domain; N_{replay} is the size of the replay memory; the target Q-network parameters θ^- are only updated every C iterations with the Q-network parameters θ and are held fixed between updates; the variable i^l corresponds to a mini-batch of tuples (s,a) taken randomly in the replay memory and the variable o^l is the corresponding target value for each tuple.

costs always underestimates the optimal cost, and therefore convergence is not assured [8, 9]. These convergence problems are verified experimentally and it has been reported that convergence may be slow or even unreliable with this online update rule [10].

We offer an analysis based on the complexity of the policy class [11] via a parallel with machine learning. High complexity machine learning techniques have the ability to represent their training set well, but at risk of overfitting. In contrast, models with lower complexity do not tend to overfit, but may fail to capture important features. In the reinforcement learning setting, an equivalence with the model complexity of the learning method may be seen in the policy class complexity. Specifically, in the context of online neural fitted reinforcement learning, the discount factor as well as the neural network architecture control the number of possible policies. Similarly to the bias-variance tradeoff observed in supervised learning, reinforcement learning also faces a trade-off between large policy class complexity and low policy class complexity.

It is already known that the optimal solution may actually be found in the set of policy with a planning horizon γ smaller than the evaluation horizon γ_{eval} specified by the problem formulation in case of only inaccurate information on the actual dynamics is available a priori (e.g. in the case where a small number of tuples is available) [11]. It is also well known that the longer the planning horizon, the greater the computational expense of computing an optimal policy [12].

In the case of neural fitted value learning, the difficulty to target a high policy class complexity is severe because targeting a high discount factor leads to propagation of errors and instabilities. Some practical ways to prevent instabilities make use of a replay memory, clipping the error term, a separate target Q-network in Equation 5 and a convolutional network architecture [1]. Using the double Q-learning algorithm also help reducing overestimations of Q-value function caused by the generalization exageration of the regression method [13].

In this paper, we investigate the tradeoff between performance of the targeted policy versus stability and speed of the learning process. We investigate the possibility to reduce instabilities during learning by working on an adaptive discount factor so as to soften the errors through learning. The goal is to target a high policy class complexity while reducing error propagations during the Deep Q learning iterations.

4 Experiments

The deep Q-learning algorithm described in [1] is used as a benchmark. All hyperparameters are kept identical if not stated differently. The main modification comes from the discount factor which is increased at every epoch (250 000 steps) with the following formula:

$$\gamma_{k+1} = 1 - 0.98(1 - \gamma_k) \tag{6}$$

The learned policies are then evaluated for 125000 steps with an ϵ -greedy policy identical to the original benchmark with $\epsilon_{test}=0.05$. The reported scores are the highest average episode score of each simulation where these evaluation episodes were not truncated at 5 min. Each game is simulated 5 times for each configuration with varying seeds and results are reported in Figure 2. It can be observed that by simply using an increasing discount factor, learning is faster for four out of the five tested games and similar for the remaining game. We conclude that by starting with a low discount factor, we obtain faster policy improvement thanks to less instability. In addition, this also provides more robustness with respect to neural network weights initialisation.

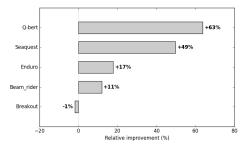


Figure 2: Summary of the results for an increasing discount factor. Reported scores are the relative improvements after 20M steps between an increasing discount factor and a constant discount factor set to its final value $\gamma=0.99$. Details are given in Appendix - Table 1.

4.1 Convergence of the neural network

We now discuss the stability of DQN. We use the experimental rule from Equation 6 and either let γ increase or keep it constant when it attains 0.99 (at 20M steps). This is illustrated on Figure 3. It is shown that increasing γ without additional care degrades severely the score obtained beyond $\gamma \approx 0.99$. By looking at the average V value, it can be seen that overestimation is particularly severe which causes the poor policies.

4.2 Further improvement with an adaptive learning rate

Since instabilities are only severe when the discount factor is high, we now study the possibility to use a more aggressive learning rate in the neural network when working at a low discount factor because potential errors would have less impact at this stage. The learning rate is then reduced along with the increasing discount factor so as to end up with a stable neural Q-learning function.

We start with a learning rate of 0.005, i.e. twice as big as the one considered in the original benchmark and we use the following simple rule at every epoch:

$$\alpha_{k+1} = 0.98\alpha_k$$

With γ that follows Equation 6 and is kept constant when it attains 0.99, we manage to improve further the score obtained on all of the games tested. Results are reported in Figure 4 and an illustration is given for two different games in Figure 5. It can be noted that the value function V decreases when γ is hold fixed and when the learning rate is lowered which is a sign of a decrease of the overestimations of the Q-value function.

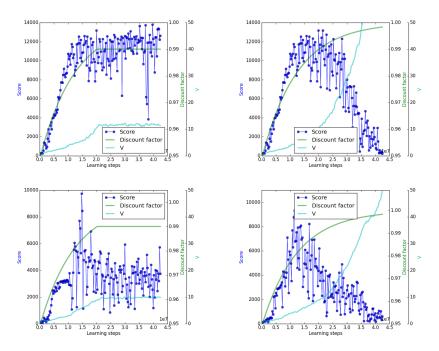


Figure 3: Illustration for the game q-bert (top) and seaquest (bottom) for a discount factor γ kept at 0.99 after 20M learning steps (on the left) as well as a discount factor that keeps increasing to a value close to 1.

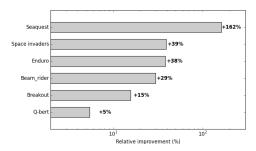


Figure 4: Summary of the results for a decreasing learning rate. Reported scores are the improvement after 50M steps when using both a dynamic discount factor and a dynamic learning rate. Details are given in Appendix - Table 2.

4.3 Exploration / Exploitation Dilemma

Errors in the policy may also have positive impacts since it increases exploration [14]. When using a lower discount factor, it actually decreases exploration and opens up the risk of falling in a local optimum in the value iteration learning. This is observed as the agent gets repeatedly a score lower than the optimal while being unable to discover some parts of the state space.

In this case, an actor-critic-type algorithm that increases the level of exploration may allow to overcome this problem. We believe that an actor critic agent that also manages adaptively the level of exploration is important to further improve deep reinforcement learning algorithms. In order to illustrate this, a simple rule has been applied in the case of the game seaquest as can be seen on Figure 6. This rule adapts the exploration during the training process in the ϵ -greedy action selection until the agent was able to get out of the local optimum.

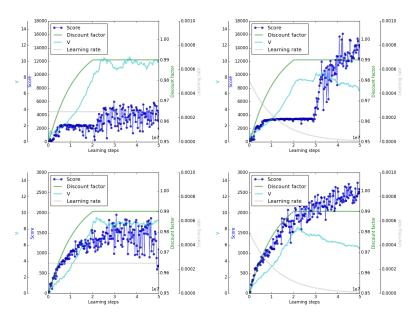


Figure 5: Illustration for the game seaquest (top) and space invaders (bottom). On the left, the deep Q-network with original parameters ($\alpha = 0.00025$) and on the right with a decreasing learning rate.

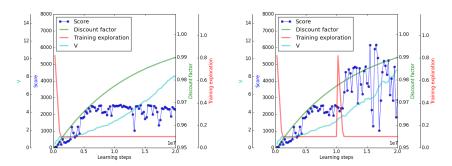


Figure 6: Illustration for the game seaquest for a case where the flat exploration rate fails to get out of a local optimum on the left. And on the right illustration where a simple rule that increases exploration allows to get out of the local optimum.

4.4 Towards an actor-critic algorithm

Following the ideas discussed above, Figure 7 represents the general update scheme that we propose to further improve the performance of deep reinforcement learning algorithms.

5 Conclusion

This paper introduced an approach to speed-up the convergence and improve the quality of the learned Q-function in deep reinforcement learning algorithms. It works by adapting the discount factor and the learning rate along the way to convergence. We used the deep Q-learning algorithms for Atari 2600 computer games as a benchmark and our approach showed improved performances for the 6 tested games. These results motivate further experiments, in particular it would be interesting to develop an automatic way to adapt online the discount factor along with the learning rate and possibly the level of exploration. It would also be of interest to combine this approach with the recent advances in deep reinforcement learning: the massively parallel architecture "Gorilla" [15],

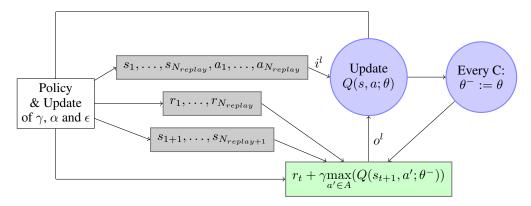


Figure 7: Sketch of the actor-critic-type algorithm for deep reinforcement learning that manages adaptively the discount rate γ , the learning rate α as well as the level of exploration ϵ .

the double Q-learning algorithm [13], the prioritized experience replay [16], the dueling network architecture [17] or a recurrent architecture such as [18].

References

- [1] 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.
- [2] Richard S Sutton and Andrew G Barto. *Introduction to reinforcement learning*. MIT Press, 1998.
- [3] Lucian Busoniu, Robert Babuska, Bart De Schutter, and Damien Ernst. *Reinforcement learning and dynamic programming using function approximators*, volume 39. CRC press, 2010.
- [4] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *arXiv preprint arXiv:1207.4708*, 2012.
- [5] Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In *ICML*, pages 30–37, 1995.
- [6] Walter Mischel, Ebbe B Ebbesen, and Antonette Raskoff Zeiss. Cognitive and attentional mechanisms in delay of gratification. *Journal of personality and social psychology*, 21(2):204, 1972.
- [7] Richard Ernest Bellman and Stuart E Dreyfus. Applied dynamic programming. 1962.
- [8] John N Tsitsiklis and Benjamin Van Roy. An analysis of temporal-difference learning with function approximation. *Automatic Control, IEEE Transactions on*, 42(5):674–690, 1997.
- [9] Geoffrey J Gordon. Approximate solutions to markov decision processes. *Robotics Institute*, page 228, 1999.
- [10] Martin Riedmiller. Neural fitted Q iteration–first experiences with a data efficient neural reinforcement learning method. In *Machine Learning: ECML 2005*, pages 317–328. Springer, 2005.
- [11] Nan Jiang, Alex Kulesza, Satinder Singh, and Richard Lewis. The dependence of effective planning horizon on model accuracy. In *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*, pages 1181–1189. International Foundation for Autonomous Agents and Multiagent Systems, 2015.
- [12] Michael Kearns, Yishay Mansour, and Andrew Y Ng. A sparse sampling algorithm for near-optimal planning in large markov decision processes. *Machine Learning*, 49(2-3):193–208, 2002.
- [13] Hado Van Hasselt, Arthur Guez, and Silver David. Deep reinforcement learning with double Q-learning. *arXiv preprint arXiv:1509.06461*, 2015.

- [14] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. *Journal of artificial intelligence research*, pages 237–285, 1996.
- [15] Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, et al. Massively parallel methods for deep reinforcement learning. arXiv preprint arXiv:1507.04296, 2015.
- [16] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.
- [17] Ziyu Wang, Nando de Freitas, and Marc Lanctot. Dueling network architectures for deep reinforcement learning. *arXiv preprint arXiv:1511.06581*, 2015.
- [18] Matthew Hausknecht and Peter Stone. Deep recurrent Q-learning for partially observable MDPs. *arXiv preprint arXiv:1507.06527*, 2015.

A Appendix

	Seaquest	Q-bert	Breakout	Enduro	Beam rider
$\gamma = 0.99, 20 \text{M steps}$	4766	7930	346	817	8379
	4774	8262	359	950	9013
	2941	7527	372	821	9201
	3597	8214	374	863	8496
	4280	7867	365	777	8408
Āverage	4072	7960	363	846	8699
Increasing γ , 20M steps	2570	13073	351	929	9011
	2575	12873	361	1031	10160
	11717	13351	362	1099	9880
	11030	12828	354	978	9263
	2583	13063	351	945	10389
Average	6095	13031	356	996	9741

Table 1: Summary of the results for an increasing discount factor.

	Seaquest	Q-bert	Breakout	Enduro	Beam rider	Space invaders
Fixed γ and α , 50M steps	6150	12958	401	817	9606	1976
Varying γ , fixed α , 50M steps	6016	13997	389	929	10677	1954
Variable γ and α , 50M steps	16124	14996	423	1129	12473	2750

Table 2: Summary of the results for an increasing discount factor associated with a decreasing learning rate.