



**HAL**  
open science

# Performance Bounds for Lambda Policy Iteration and Application to the Game of Tetris

Bruno Scherrer

► **To cite this version:**

Bruno Scherrer. Performance Bounds for Lambda Policy Iteration and Application to the Game of Tetris. [Research Report] RR-6348, 2011, pp.54. inria-00185271v4

**HAL Id: inria-00185271**

**<https://inria.hal.science/inria-00185271v4>**

Submitted on 3 Oct 2011 (v4), last revised 11 Oct 2011 (v5)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Performance Bounds for $\lambda$ Policy Iteration and Application to the Game of Tetris

**Bruno Scherrer**

*Maia Team, INRIA Lorraine*

*615 rue du Jardin Botanique*

*54600 Villers-les-Nancy*

*FRANCE*

BRUNO.SCHERRER@INRIA.FR

## Abstract

We consider the discrete-time infinite-horizon optimal control problem formalized by Markov Decision Processes (Puterman, 1994; Bertsekas and Tsitsiklis, 1996). We revisit the work of Bertsekas and Ioffe (1996), that introduced  $\lambda$  Policy Iteration, a family of algorithms parameterized by  $\lambda$  that generalizes the standard algorithms Value Iteration and Policy Iteration, and has some deep connections with TD( $\lambda$ ) introduced by Sutton and Barto (1998). We deepen the original theory developed by the authors by providing convergence rate bounds which generalize standard bounds for Value Iteration described for instance by Puterman (1994). Then, the main contribution of this paper is to develop the theory of this algorithm when it is used in an approximate form and show that this is sound. Doing so, we extend and unify the separate analyses developed by Munos for Approximate Value Iteration (Munos, 2007) and Approximate Policy Iteration (Munos, 2003). Eventually, we revisit the use of this algorithm in the training of a Tetris playing controller as originally done by Bertsekas and Ioffe (1996). We provide an original performance bound that can be applied to such an undiscounted control problem. Our empirical results are different from those of Bertsekas and Ioffe (which were originally qualified as “paradoxical” and “intriguing”), and much more conform to what one would expect from a learning experiment. We discuss what we believe are the reasons for such a difference.

**Keywords:** Stochastic Optimal Control, Reinforcement Learning, Markov Decision Processes, Analysis of Algorithms, Performance Bounds.

## 1. Introduction

We consider the discrete-time infinite-horizon optimal control problem formalized by Markov Decision Processes (Puterman, 1994; Bertsekas and Tsitsiklis, 1996). We revisit the  $\lambda$  Policy Iteration algorithm introduced by Bertsekas and Ioffe (1996) (also published in the reference textbook of Bertsekas and Tsitsiklis (1996)<sup>1</sup>), that (as the authors then stated) “*is primarily motivated by the case of large and complex problems where the use of approximation is essential*”. It is a family of algorithms parameterized by  $\lambda$  that generalizes the standard Dynamic Programming algorithms Value Iteration (which corresponds to the case  $\lambda = 0$ ) and Policy Iteration (case  $\lambda = 1$ ), and has some deep connections with TD( $\lambda$ ) algorithms

---

1. The reference (Bertsekas and Ioffe, 1996) being historically anterior to (Bertsekas and Tsitsiklis, 1996), we will only refer to the former in the rest of the paper.

that are well known to the Reinforcement Learning community (Sutton and Barto, 1998; Bertsekas and Tsitsiklis, 1996).

In their original paper, Bertsekas and Ioffe (1996) show the convergence of  $\lambda$  Policy Iteration when it is run without error and provide its *asymptotic* convergence rate. The authors also describe a case study involving an instance of Approximate  $\lambda$  Policy Iteration, but neither their paper nor (to the best of our knowledge) any subsequent work studies the theoretical soundness of doing so. In this paper, we extend the theory on this algorithm in several ways. We derive its *non-asymptotic* convergence rate when it is run without error. More importantly, we develop the theory of  $\lambda$  Policy Iteration for its main purpose, that is (cf. the above quote) when it is run in an approximate form, and prove that such an approach is sound: we show that the loss of using the greedy policy with respect to the current value estimate can be made arbitrarily small by controlling the error made during the iterations.

The rest of the paper is organized as follows. In Section 2, we introduce the framework of Markov Decision Processes and describe the two standard and widely known algorithms, Value Iteration and Policy Iteration, along with some of their state-of-the-art analysis in exact and approximate form. Section 3 introduces  $\lambda$  Policy Iteration in an original way that makes its connection with Value Iteration and Policy Iteration obvious, and discusses its close connection with Reinforcement Learning. We recall the main results obtained by Bertsekas and Ioffe (1996) (convergence and asymptotic rate of convergence of the exact algorithm). At this point of the paper, we naturally describe how one expects that the properties of Value Iteration ( $\lambda = 0$ ) and Policy Iteration ( $\lambda = 1$ ) described in Section 2 may translate for general  $\lambda$ . The precise statements of our results are the topic of the next two Sections: Section 5 contains our new results on Exact  $\lambda$  Policy Iteration and Section 6 those on Approximate  $\lambda$  Policy Iteration<sup>2</sup>. Last but not least, Section 7 revisits the empirical part of the work of Bertsekas and Ioffe (1996), where an approximate version of  $\lambda$  Policy Iteration is used for training a Tetris controller.

**General Notations** The analysis we will describe in this article relies on a few notations, such as several norms and seminorms, that we need to define precisely before we can go further. Let  $X$  be a finite space. Let  $u$  denote a real-valued function on  $X$ , which can be seen as a vector of dimension  $|X|$ . Let  $e$  denote the vector of which all components are 1.  $\mu$  denotes a distribution on  $X$ . We consider the **weighted  $L_p$  norm**:

$$\|u\|_{p,\mu} := \left( \sum_x \mu(x) |u(x)|^p \right)^{1/p} = (\mu^T |u|^p)^{1/p}$$

where  $|u|^p$  denotes the componentwise absolute value and exponentiation of  $u$ . We will write  $\|\cdot\|_p$  the **unweighted  $L_p$  norm** (i.e. with uniform distribution  $\mu$ ).  $\|\cdot\|_\infty$  is the max-norm:

$$\|u\|_\infty := \max_x |u(x)| = \lim_{p \rightarrow \infty} \|u\|_p.$$

We write  $\text{span}_\infty[\cdot]$  the **span seminorm** (as for instance defined in (Puterman, 1994)):

$$\text{span}_\infty[u] := \max_x u(x) - \min_x u(x).$$

---

2. Section 6 is probably the place where the reader familiar with Approximate Dynamic Programming would quickly want to jump.

It can be seen that

$$\text{span}_\infty [u] = 2 \min_a \|u - ae\|_\infty.$$

We propose to generalize the span seminorm definition for any  $p$  as follows:

$$\text{span}_{p,\mu} [u] := 2 \min_a \|u - ae\|_{p,\mu}$$

It is clear that it is a seminorm (i.e. it is non-negative, it satisfies the triangle inequality and  $\text{span}_*[au] = |a|\text{span}_*[u]$ ). It is not a norm because it is zero for all constant functions.

The error bounds we will derive in this paper are expressed in terms of some span seminorm. The following relations

$$\begin{cases} \text{span}_p [u] & \leq 2 \|u\|_p & \leq 2 \|u\|_\infty \\ \text{span}_{p,\mu} [u] & \leq 2 \|u\|_{p,\mu} & \leq 2 \|u\|_\infty \\ \text{span}_\infty [u] & \leq 2 \|u\|_\infty \end{cases} \quad (1)$$

show how to deduce error bounds involving the (more standard)  $L_p$  and max norms. Since the span seminorm can be zero for non zero (constant) vectors, there is no relation that would enable us to derive error bounds in span seminorm from a  $L_p$  or a max norm. Bounding an error with the span seminorm is in this sense stronger and this constitutes our motivation for using it.

## 2. Framework and Standard Algorithms

In this section, we begin by providing a short description of the framework of Markov Decision Processes we consider throughout the paper. We go on by describing the two main algorithms, Value Iteration and Policy Iteration, for solving the related problem.

### 2.1 Markov Decision Processes

We consider a discrete-time dynamic system whose state transition depends on a control. We assume that there is a **state space**  $X$  of finite size  $N$ . When at state  $i \in \{1, \dots, N\}$ , the control is chosen from a finite **control space**  $A$ . The control  $a \in A$  specifies the **transition probability**  $p_{ij}(a)$  to the next state  $j$ . At the  $k$ th iteration, the system is given a reward  $\gamma^k r(i, a, j)$  where  $r$  is the instantaneous **reward function**, and  $0 < \gamma < 1$  is a discount factor. The tuple  $\langle X, A, p, r, \gamma \rangle$  is known as a **Markov Decision Process (MDP)** (Puterman, 1994; Bertsekas and Tsitsiklis, 1996).

We are interested in stationary deterministic policies, that is functions  $\pi : X \rightarrow A$  which map states into controls<sup>3</sup>. Writing  $i_k$  the state at time  $k$ , the **value of policy**  $\pi$  at state  $i$  is defined as the total expected return while following a policy  $\pi$  from  $i$ , that is

$$v^\pi(i) := \lim_{N \rightarrow \infty} E_\pi \left[ \sum_{k=0}^{N-1} \gamma^k r(i_k, \pi(i_k), i_{k+1}) \mid i_0 = i \right] \quad (2)$$

---

3. Restricting our attention to stationary deterministic policies is not a limitation. Indeed, for the optimality criterion to be defined soon, it can be shown that there exists at least one stationary deterministic policy which is optimal (Puterman, 1994).

where  $E_\pi$  denotes the expectation conditional on the fact that the actions are selected with the policy  $\pi$  (that is, for all  $k$ ,  $i_{k+1}$  is reached from  $i_k$  with probability  $p_{i_k i_{k+1}}(\pi(i_k))$ ). The **optimal value** starting from state  $i$  is defined as

$$v_*(i) := \max_{\pi} v^\pi(i).$$

We write  $P^\pi$  the  $N \times N$  stochastic matrix whose elements are  $p_{ij}(\pi(i))$  and  $r^\pi$  the vector whose components are  $\sum_j p_{ij}(\pi(i))r(i, \pi(i), j)$ .  $v^\pi$  and  $v_*$  can be seen as vectors on  $X$ . It is well-known that  $v^\pi$  solves the following Bellman Equation:

$$v^\pi = r^\pi + \gamma P^\pi v^\pi.$$

$v^\pi$  is a fixed point of the linear operator  $\mathcal{T}^\pi v := r^\pi + \gamma P^\pi v$ . As  $P^\pi$  is a stochastic matrix, its eigenvalues cannot be greater than 1, and consequently  $I - \gamma P^\pi$  is invertible. This implies that

$$v^\pi = (I - \gamma P^\pi)^{-1} r^\pi = \sum_{i=0}^{\infty} (\gamma P^\pi)^i r^\pi. \quad (3)$$

It is also well-known that  $v_*$  satisfies the following Bellman Equation:

$$v_* = \max_{\pi} (r^\pi + \gamma P^\pi v_*) = \max_{\pi} \mathcal{T}^\pi v_*$$

where the max operator is componentwise. In other words,  $v_*$  is a fixed point of the nonlinear operator  $\mathcal{T} v := \max_{\pi} \mathcal{T}^\pi v$ . For any value vector  $v$ , we call a **greedy policy with respect to the value**  $v$  a policy  $\pi$  that satisfies:

$$\pi \in \arg \max_{\pi'} \mathcal{T}^{\pi'} v$$

or equivalently  $\mathcal{T}^\pi v = \mathcal{T} v$ . We will write, with some abuse of notation<sup>4</sup>  $\text{greedy}(v)$  any policy that is greedy with respect to  $v$ . The notions of optimal value function and greedy policies are fundamental to optimal control because of the following property: any policy  $\pi_*$  that is greedy with respect to the optimal value is an **optimal policy** and its value  $v^{\pi_*}$  equals  $v_*$ .

It is well-known that the operators  $\mathcal{T}^\pi$  and  $\mathcal{T}$  are  $\gamma$ -contraction mappings with respect to the max-norm. In what follows we only write what this means for the Bellman operator  $\mathcal{T}$  but the same holds for  $\mathcal{T}^\pi$ . Being a  $\gamma$ -contraction mapping for the max-norm means that for all pairs of vectors  $(v, w)$ ,

$$\|\mathcal{T}v - \mathcal{T}w\|_\infty \leq \gamma \|v - w\|_\infty.$$

This ensures that the fixed point  $v_*$  of  $\mathcal{T}$  exists and is unique. Furthermore, for any initial vector  $v_0$ ,

$$\lim_{k \rightarrow \infty} (\mathcal{T})^k v_0 = v_*. \quad (4)$$

Given an MDP, standard algorithmic solutions for computing an optimal value/policy (which dates back to the 1950s, see for instance Puterman (1994, and the references therein)) are Value Iteration and Policy Iteration. The rest of this section describes both of these algorithms with some of the relevant properties for the subject of this paper.

---

4. There might be several policies that are greedy with respect to some value  $v$ .

## 2.2 Value Iteration

The Value Iteration algorithms for computing the value of a policy  $\pi$  and the value of the optimal policy  $\pi_*$  rely on Equation 4. Algorithm 1 provides a description of Value Iteration for computing an optimal policy (replace  $\mathcal{T}$  by  $\mathcal{T}^\pi$  in it and one gets Value Iteration for computing the value of some policy  $\pi$ ). In this description, we have introduced a term  $\epsilon_k$

---

### Algorithm 1 Value Iteration

---

**Input:** An MDP, an initial value  $v_0$

**Output:** An (approximately) optimal policy

$k \leftarrow 0$

**repeat**

$v_{k+1} \leftarrow \mathcal{T}v_k + \epsilon_{k+1}$  // Update the value

$k \leftarrow k + 1$

**until** some stopping criterion

**Return** greedy( $v_k$ )

---

which stands for several possible sources of error at each iteration: this error might be the computer round off, the fact that we use an approximate architecture for representing  $v$ , a stochastic approximation of  $P^{\pi_k}$ , etc... or a combination of these. In what follows, when we talk about the *Exact* version of an algorithm, this means that  $\epsilon_k = 0$  for all  $k$ .

**Properties of Exact Value Iteration** It is well-known that the contraction property induces some interesting properties for Exact Value Iteration. We have already mentioned that contraction implies the asymptotic convergence (Equation 4). It can also be inferred that there is at least a linear rate of convergence: for all reference iteration  $k_0$ , and for all  $k \geq k_0$ ,

$$\|v_* - v_k\|_\infty \leq \gamma^{k-k_0} \|v_* - v_{k_0}\|_\infty.$$

Even more interestingly, it is possible to derive a performance bound, that is a bound of the difference between the real value of a policy produced by the algorithm and the value of the optimal policy  $\pi_*$  Puterman (1994). Let  $\pi_k$  denote the policy that is greedy with respect to  $v_{k-1}$ . Then, for all reference iteration  $k_0$ , and for all  $k \geq k_0$ ,

$$\|v_* - v^{\pi_k}\|_\infty \leq \frac{2\gamma^{k-k_0}}{1-\gamma} \|\mathcal{T}v_{k_0} - v_{k_0}\|_\infty = \frac{2\gamma^{k-k_0}}{1-\gamma} \|v_{k_0+1} - v_{k_0}\|_\infty.$$

This fact is of considerable importance computationally since it provides a stopping criterion: taking  $k = k_0 + 1$ , we see that if  $\|v_{k_0+1} - v_{k_0}\|_\infty \leq \frac{1-\gamma}{2\gamma}\epsilon$ , then  $\|v_* - v^{\pi_{k_0+1}}\|_\infty \leq \epsilon$ .

It is somewhat less known that the Bellman operators  $\mathcal{T}$  and  $\mathcal{T}^\pi$  are also contraction mapping with respect to the  $\text{span}_\infty$  seminorm (Puterman, 1994). This means that there exists a variant of the above Equation involving the span seminorm instead of the max-norm. For instance, such a fact provides the following stopping criterion:

**Proposition 1 (Stopping Cond. for Exact Value Iteration (Puterman, 1994))**

*If at some iteration  $k_0$ , the difference between two subsequent iterations satisfies*

$$\text{span}_\infty [v_{k_0+1} - v_{k_0}] \leq \frac{1-\gamma}{\gamma}\epsilon,$$

then the greedy policy  $\pi_{k_0+1}$  with respect to  $v_{k_0}$  is  $\epsilon$ -optimal:  $\|v_* - v^{\pi_{k_0+1}}\|_\infty \leq \epsilon$ .

This latter stopping criterion is finer since, from the relation between the span seminorm and the norm (Equation 1) it implies the former.

**Properties of Approximate Value Iteration (AVI)** When considering large Markov Decision Processes, one cannot usually implement an exact version of Value Iteration. In such a case  $\epsilon_k \neq 0$ . In general, the algorithm does not converge anymore but it is possible to study its asymptotic behaviour. The most well-known result is due to Bertsekas and Tsitsiklis (1996): If the approximation errors are uniformly bounded, the difference between the asymptotic performance of policies  $\pi_{k+1}$  greedy with respect to  $v_k$  satisfies

$$\limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty \leq \frac{2\gamma}{(1-\gamma)^2} \sup_{k \geq 0} \|\epsilon_k\|_\infty. \quad (5)$$

Munos (2003, 2007) has recently argued that, since most supervised learning algorithms (such as least square regression) that are used in practice for approximating each iterate of Value Iteration control some  $L_p$  norm, it would be more interesting to have an analogue of the above result where the approximation error  $\epsilon_k$  is expressed in terms of the  $L_p$  norm. Munos (2007) actually showed how to do this. The idea is to analyze the *componentwise* asymptotic behaviour of Approximate Value Iteration, from which it is possible to derive  $L_p$  bounds for any  $p$ . Write  $P_k = P^{\pi_k}$  the stochastic matrix corresponding to the policy  $\pi_k$  which is greedy with respect to  $v_{k-1}$ ,  $P_*$  the stochastic matrix corresponding to the (unknown) optimal policy  $\pi_*$ . Munos (2007) showed the following Lemma:

**Lemma 2 (Asymptotic Componentwise Performance of AVI (Munos, 2007))**

*The following matrices*

$$\begin{aligned} Q_{kj} &:= (1-\gamma)(I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_{j+1} \\ Q'_{kj} &:= (1-\gamma)(I - \gamma P_k)^{-1} (P_*)^{k-j} \end{aligned}$$

*are stochastic and the asymptotic performance of the policies generated by Approximate Value Iteration satisfies*

$$\limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty \leq \limsup_{k \rightarrow \infty} \frac{1}{1-\gamma} \sum_{j=0}^{k-1} \gamma^{k-j} [Q_{kj} - Q'_{kj}] \epsilon_j.$$

From the above componentwise bound, it is possible<sup>5</sup> to derive the following  $L_p$  bounds.

**Proposition 3 (Asymptotic Performance of AVI (1/2))**

*Choose any  $p$  and any distribution  $\mu$ . Consider the notations of Lemma 2. Then for all  $k \geq j \geq 0$ ,*

$$\mu_{kj} := \frac{1}{2} (Q_{jk} + Q'_{jk})^T \mu$$

---

5. This result is not explicitly stated by Munos (2007), but using a technique of another of his articles (Munos, 2003), it is straightforward to derive from Lemma 2. The current paper will anyway generalize this result (in Proposition 24 page 23).

are distributions and the asymptotic performance of policies generated by Value Iteration satisfies

$$\limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{p,\mu} \leq \frac{2\gamma}{(1-\gamma)^2} \sup_{k \geq j \geq 0} \|\epsilon_k\|_{p,\mu_{kj}}$$

As the above bounds rely on the partially unknown matrices  $Q_{kj}$  and  $Q'_{kj}$ , Munos (2003, 2007) introduced some assumption in terms of **concentration coefficient** to remove this dependency. Assume there exists a distribution  $\nu$  and a real number  $C(\nu)$  such that

$$C(\nu) := \max_{i,j,a} \frac{p_{ij}(a)}{\nu(j)}. \quad (6)$$

For instance, if one chooses the uniform law  $\nu$ , then there always exists such a  $C(\nu) \in (1, N)$  where  $N$  is the size of the state space (see (Munos, 2003, 2007) for more discussion on this coefficient). This allows to derive the following performance bounds on the max-norm of the loss.

**Proposition 4 (Asymptotic Performance of AVI (2/2) (Munos, 2007))**

Let  $C(\nu)$  be the concentration coefficient defined in Equation 6. The asymptotic performance of the policies generated by Approximate Value Iteration satisfies

$$\limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{\infty} \leq \frac{2\gamma (C(\nu))^{1/p}}{(1-\gamma)^2} \sup_{k \geq 0} \|\epsilon_k\|_{p,\nu}.$$

The main difference between the bounds of Propositions 3 and 4 and that of Bertsekas and Tsitsiklis (Equation 5) is that the approximation error  $\epsilon_k$  is measured by the weighted  $L_p$  norm. As  $\lim_{p \rightarrow \infty} \|\cdot\|_{p,\mu} \leq \|\cdot\|_{\infty}$  and  $(C(\nu))^{1/p} \xrightarrow{p \rightarrow \infty} 1$ , Munos's results are strictly finer.

There is in general no guarantee that AVI converges. AVI is known to converge for specific approximation architectures known as averagers (Gordon, 1995) which include state aggregation (see (Van Roy, 2006) for a very fine approximation bound in this specific case). Also, convergence may just occur experimentally. Suppose the sequence  $(v_k)_{k \geq 0}$  tends to some  $v$ . Write  $\pi$  the corresponding greedy policy. Notice this implies that  $(\epsilon_k)_{k \geq 0}$  tends to  $v - \mathcal{T}v$ , which is known as the **Bellman residual**. The above bounds can be improved by a factor  $\frac{1}{1-\gamma}$ . It is indeed known (e.g. Williams and Baird (1993)) that

$$\|v_* - v^{\pi}\|_{\infty} \leq \frac{2\gamma}{1-\gamma} \|v - \mathcal{T}v\|_{\infty}$$

and, with the same notations as above, Munos (2007) derived the analogous finer  $L_p$  bound:

**Corollary 5 (Performance of AVI in case of convergence (Munos, 2007))**

Let  $C(\nu)$  be the concentration coefficient defined in Equation 6. Suppose  $(v_k)$  tends to some  $v$ . Write  $\pi$  the corresponding greedy policy. Then

$$\|v_* - v^{\pi}\|_{\infty} \leq \frac{2\gamma (C(\nu))^{1/p}}{1-\gamma} \|v - \mathcal{T}v\|_{p,\nu}.$$

Eventually, let us mention that Munos (2007) and Farahmand *et al.* (2010) consider some *finer* performance bounds (in weighted  $L_p$  norm) using some *finer* concentration coefficients. We won't discuss them in this paper and we recommend the interested reader to go through these references for further details.



### 2.3 Policy Iteration

Policy Iteration is an alternative method for computing an optimal policy for an infinite-horizon discounted Markov Decision Process. This algorithm is based on the following

---

**Algorithm 2** Policy Iteration
 

---

**Input:** An MDP, an initial policy  $\pi_0$

**Output:** An (approximately) optimal policy

$k \leftarrow 0$

**repeat**

$v_k \leftarrow (I - \gamma P^{\pi_k})^{-1} r^{\pi_k} + \epsilon_k$  // Estimate the value of  $\pi_k$

$\pi_{k+1} \leftarrow \text{greedy}(v_k)$  // Update the policy

$k \leftarrow k + 1$

**until** some stopping criterion

**Return**  $\pi_k$

---

property: if  $\pi$  is some policy, then any policy  $\pi'$  that is greedy with respect to the value of  $\pi$ , i.e. any  $\pi'$  satisfying  $\pi' = \text{greedy}(v^\pi)$ , is better than  $\pi$  in the sense that  $v^{\pi'} \geq v^\pi$ . Policy Iteration exploits this property in order to generate a sequence of policies with increasing values. It is described in Algorithm 2. Note that we use the analytical form of the value of a policy given by Equation 3. Also, as for Value Iteration, our description includes a potential error  $\epsilon_k$  term each time the value of a policy is estimated.

**Properties of Exact Policy Iteration** When the state space and the control spaces are finite, it is known that Exact Policy Iteration converges to an optimal policy  $\pi_*$  in a finite number of iterations. In infinite state spaces, if the function  $v \mapsto P^{\text{greedy}(v)}$  is Lipschitz, then it can be shown that Policy Iteration has a quadratic convergence rate (Puterman, 1994). However, to our knowledge, and contrary to Value Iteration, finite-time stopping conditions such as that of Proposition 1 are not widely known for Policy Iteration (though they appear implicitly in some recent works on Approximate Policy Iteration (Antos *et al.*, 2007, 2008; Farahmand *et al.*, 2010; Lazaric *et al.*, 2010)).

**Properties of Approximate Policy Iteration (API)** For problems of interest, one usually uses Policy Iteration in an approximate form, i.e. with  $\epsilon_k \neq 0$ . Results similar to those we presented for Approximate Value Iteration exist for Approximate Policy Iteration. As soon as there is some error  $\epsilon_k \neq 0$ , the algorithm does not necessarily converge anymore but there is an analogue of Equation 5 which is also due to Bertsekas and Tsitsiklis (1996): If the approximation errors are uniformly bounded, then the difference between the asymptotic performance of policies  $\pi_{k+1}$  greedy with respect to  $v_k$  and the optimal policy is

$$\limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty \leq \frac{2\gamma}{(1-\gamma)^2} \sup_{k \geq 0} \|\epsilon_k\|_\infty. \quad (7)$$

As for Value Iteration, Munos has extended this result so that one can get bounds involving the  $L_p$  norm. He also showed how to relate the performance analysis to the Bellman residual  $v_k - \mathcal{T}^{\pi_k} v_k$  that says how much  $v_k$  approximates the real value of the policy  $\pi_k$ ; this is for instance interesting when the evaluation step of Approximate Policy Iteration involves the

minimization of the norm of this Bellman residual (see (Munos, 2003)). It is important to note that this Bellman residual is different from the one we introduced in the previous section (we then considered  $v_k - \mathcal{T}v_k = v_k - \mathcal{T}^{\pi_{k+1}}v_k$  where  $\pi_{k+1}$  is greedy with respect to  $v_k$ ). To avoid confusion, and because it is related to some specific policy, we will call  $v_k - \mathcal{T}^{\pi_k}v_k$  the **Policy Bellman residual**. Munos started by deriving a componentwise analysis. Write  $P_k = P^{\pi_k}$  the stochastic matrix corresponding to the policy  $\pi_k$  which is greedy with respect to  $v_{k-1}$ ,  $P_*$  the stochastic matrix corresponding to the (unknown) optimal policy  $\pi_*$ .

**Lemma 6 (Asymptotic Componentwise Performance of API (Munos, 2003))**

*The following matrices*

$$\begin{aligned} R_k &:= (1 - \gamma)^2 (I - \gamma P_*)^{-1} P_{k+1} (I - \gamma P_{k+1})^{-1} \\ R'_k &:= (1 - \gamma)^2 (I - \gamma P_*)^{-1} [P_* + \gamma P_{k+1} (I - \gamma P_{k+1})^{-1} P_k] \\ R''_k &:= (1 - \gamma)^2 (I - \gamma P_*)^{-1} P_* (I - \gamma P_k) \end{aligned}$$

*are stochastic and the asymptotic performance of the policies generated by Approximate Policy Iteration satisfies*

$$\begin{aligned} \limsup_{k \rightarrow \infty} v_* - v^{\pi_k} &\leq \frac{2\gamma}{(1 - \gamma)^2} \limsup_{k \rightarrow \infty} [R_k - R'_k] \epsilon_k \\ \limsup_{k \rightarrow \infty} v_* - v^{\pi_k} &\leq \frac{2\gamma}{(1 - \gamma)^2} \limsup_{k \rightarrow \infty} [R_k - R''_k] (v_k - \mathcal{T}^{\pi_k}v_k). \end{aligned}$$

As for Value Iteration, the above componentwise bound leads to the following  $L_p$  bounds.

**Proposition 7 (Asymptotic Performance of API (1/2) (Munos, 2003))**

*Choose any  $p$  and any distribution  $\mu$ . Consider the notations of Lemma 6. For all  $k \geq 0$ ,*

$$\begin{aligned} \mu_k &= \frac{1}{2} (R_k + R'_k)^T \mu \quad \text{and} \\ \mu'_k &= \frac{1}{2} (R_k + R''_k)^T \mu \end{aligned}$$

*are distributions and the asymptotic performance of the policies generated by Approximate Policy Iteration satisfies*

$$\begin{aligned} \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{p, \mu} &\leq \frac{2\gamma}{(1 - \gamma)^2} \limsup_{k \rightarrow \infty} \|\epsilon_k\|_{p, \mu_k} \\ \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{p, \mu} &\leq \frac{2\gamma}{(1 - \gamma)^2} \limsup_{k \rightarrow \infty} \|v_k - \mathcal{T}^{\pi_k}v_k\|_{p, \mu'_k}. \end{aligned}$$

Using the concentration coefficient  $C(\nu)$  introduced in the previous section (Equation 6), it is also possible to show<sup>6</sup> the following  $L_\infty/L_p$  bounds:

---

6. Similarly to footnote 5, this result is not explicitly stated by Munos (2003) but using techniques of another of his articles (Munos, 2007), it is straightforward to derive from Lemma 6. The current paper will anyway generalize this result (in Proposition 26 page 24).

**Proposition 8 (Asymptotic Performance of API (2/2))**

Let  $C(\nu)$  be the concentration coefficient defined in Equation 6. The asymptotic performance of the policies generated by Approximate Policy Iteration satisfies

$$\begin{aligned} \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty &\leq \frac{2\gamma (C(\nu))^{1/p}}{(1-\gamma)^2} \limsup_{k \rightarrow \infty} \|\epsilon_k\|_{p,\nu} \\ \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty &\leq \frac{2\gamma (C(\nu))^{1/p}}{(1-\gamma)^2} \limsup_{k \rightarrow \infty} \|v_k - \mathcal{T}^{\pi_k} v_k\|_{p,\nu}. \end{aligned}$$

Again, the bounds of Propositions 7 and 8 with respect to the approximation error  $\epsilon_k$  are finer than that of Bertsekas and Tsitsiklis (Equation 7). Compared to the similar result for Approximate Value Iteration (Propositions 3 and 4) where the bound depends on a *uniform* error bound ( $\forall k, \|\epsilon_k\|_{p,\nu} \leq \epsilon$ ), the above bounds have the nice property that they only depend on *asymptotic* errors/residuals.

Finally, as for Approximate Value Iteration, a better bound (by a factor  $\frac{1}{1-\gamma}$ ) might be obtained if the sequence of policies happens to converge. It can be shown (Munos, 2003, Remark 4 page 7) that:

**Corollary 9 (Performance of API in case of convergence)**

Let  $C(\nu)$  be the concentration coefficient defined in Equation 6. If the sequence of policies  $(\pi_k)$  converges to some  $\pi$ , then

$$\begin{aligned} v_* - v^\pi &\leq \frac{2\gamma (C(\nu))^{1/p}}{1-\gamma} \limsup_{k \rightarrow \infty} \|\epsilon_k\|_{p,\nu} \\ v_* - v^\pi &\leq \frac{2\gamma (C(\nu))^{1/p}}{1-\gamma} \limsup_{k \rightarrow \infty} \|v_k - \mathcal{T}^{\pi_k} v_k\|_{p,\nu}. \end{aligned}$$

After this tour of results for Value and Policy Iteration, we will now introduce the algorithm we study in this paper.

**3.  $\lambda$  Policy Iteration**

Though all the results we have emphasized so far are strongly related (and even sometimes identical, compare e.g. Equations 5 and 7), they were surprisingly proved independently. In this section, we describe the family of algorithms “ $\lambda$  Policy Iteration”<sup>7</sup> introduced by Bertsekas and Ioffe (1996) parameterized by a coefficient  $\lambda \in (0, 1)$ , that generalizes them both. When  $\lambda = 0$ ,  $\lambda$  Policy Iteration will reduce to Value Iteration while it will reduce to Policy Iteration when  $\lambda = 1$ . We will also recall the fact discussed by Bertsekas and Ioffe (1996) that  $\lambda$  Policy Iteration draws some connections with Temporal Difference algorithms (Sutton and Barto, 1998).

**3.1 The Algorithm**

We begin by giving some intuition about how one can make a connection between Value Iteration and Policy Iteration. For simplicity, let us temporarily forget about the error term

---

7. It was also called “Temporal Difference-Based Policy Iteration” in the original paper, but we take the name  $\lambda$  Policy Iteration, as it was the name picked by most subsequent works.

$\epsilon_k$ . At first sight, Value Iteration builds a sequence of value functions and Policy Iteration a sequence of policies. In fact, both algorithms can be seen as updating a sequence of value-policy pairs. With some little rewriting — by decomposing the (nonlinear) Bellman operator  $\mathcal{T}$  into 1) the maximization step and 2) the application of the (linear) Bellman operator — it can be seen that each iterate of Value Iteration is equivalent to the two following updates:

$$\begin{cases} \pi_{k+1} \leftarrow \text{greedy}(v_k) \\ v_{k+1} \leftarrow \mathcal{T}^{\pi_{k+1}} v_k \end{cases} \Leftrightarrow \begin{cases} \pi_{k+1} \leftarrow \text{greedy}(v_k) \\ v_{k+1} \leftarrow r^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} v_k. \end{cases}$$

The left hand side of the above Equation uses the operator  $\mathcal{T}^{\pi_{k+1}}$  while the right hand side uses its definition. Similarly — by inverting in Algorithm 2 the order of 1) the estimation of the value of the current policy and 2) the update of the policy, and by using the fact that the value of policy  $\pi_{k+1}$  is the fixed point of  $\mathcal{T}^{\pi_{k+1}}$  (Equation 4) — it can be argued that every iteration of Policy Iteration does the following:

$$\begin{cases} \pi_{k+1} \leftarrow \text{greedy}(v_k) \\ v_{k+1} \leftarrow (\mathcal{T}^{\pi_{k+1}})^\infty v_k \end{cases} \Leftrightarrow \begin{cases} \pi_{k+1} \leftarrow \text{greedy}(v_k) \\ v_{k+1} \leftarrow (I - \gamma P^{\pi_{k+1}})^{-1} r^{\pi_{k+1}}. \end{cases}$$

This rewriting makes both algorithms look close to each other. Both can be seen as having an estimate  $v_k$  of the value of policy  $\pi_k$ , from which they deduce a potentially better policy  $\pi_{k+1}$ . The corresponding value  $v^{\pi_{k+1}}$  of this better policy may be regarded as a target which is tracked by the next estimate  $v_{k+1}$ . The difference is in the update that enables to go from  $v_k$  to  $v_{k+1}$ : while Policy Iteration directly *jumps to* the value of  $\pi_{k+1}$  (by applying the Bellman operator  $\mathcal{T}^{\pi_{k+1}}$  an infinite number of times), Value Iteration only *makes one step* towards it (by applying  $\mathcal{T}^{\pi_{k+1}}$  only once). From this common view of Value Iteration, it is natural to introduce the well-known Modified Policy Iteration algorithm (Puterman and Shin, 1978) which *makes  $n$  steps* at each update:

$$\begin{cases} \pi_{k+1} \leftarrow \text{greedy}(v_k) \\ v_{k+1} \leftarrow (\mathcal{T}^{\pi_{k+1}})^n v_k \end{cases} \Leftrightarrow \begin{cases} \pi_{k+1} \leftarrow \text{greedy}(v_k) \\ v_{k+1} \leftarrow [I + \dots + (\gamma P^{\pi_{k+1}})^{n-1}] r^{\pi_{k+1}} + (\gamma P^{\pi_{k+1}})^n v_k. \end{cases}$$

The above common view is actually here interesting because it also leads to a natural introduction of  $\lambda$  Policy Iteration.  $\lambda$  Policy Iteration is doing a  $\lambda$ -adjustable step towards the value of  $\pi_{k+1}$ :

$$\begin{cases} \pi_{k+1} \leftarrow \text{greedy}(v_k) \\ v_{k+1} \leftarrow (1 - \lambda) \sum_{j=0}^{\infty} \lambda^j (\mathcal{T}^{\pi_{k+1}})^{j+1} v_k \end{cases} \Leftrightarrow \begin{cases} \pi_{k+1} \leftarrow \text{greedy}(v_k) \\ v_{k+1} \leftarrow (I - \lambda \gamma P^{\pi_{k+1}})^{-1} (r^{\pi_{k+1}} + (1 - \lambda) \gamma P^{\pi_{k+1}} v_k) \end{cases}$$

**Remark 10** *The equivalence between the left and the right representation of  $\lambda$  Policy Iteration needs here to be proved. For all  $k \geq 0$  and all function  $v$ , Bertsekas and Ioffe (1996) introduce the following operator<sup>8</sup>*

$$M_k \mathbf{v} := (1 - \lambda) \mathcal{T}^{\pi_{k+1}} v_k + \lambda \mathcal{T}^{\pi_{k+1}} \mathbf{v} \quad (8)$$

$$= r^{\pi_{k+1}} + (1 - \lambda) \gamma P^{\pi_{k+1}} v_k + \lambda \gamma P^{\pi_{k+1}} \mathbf{v} \quad (9)$$

and prove that

---

8. The equivalence between Equations 8 and 9 follows trivially from the definition of  $\mathcal{T}^{\pi_{k+1}}$ .

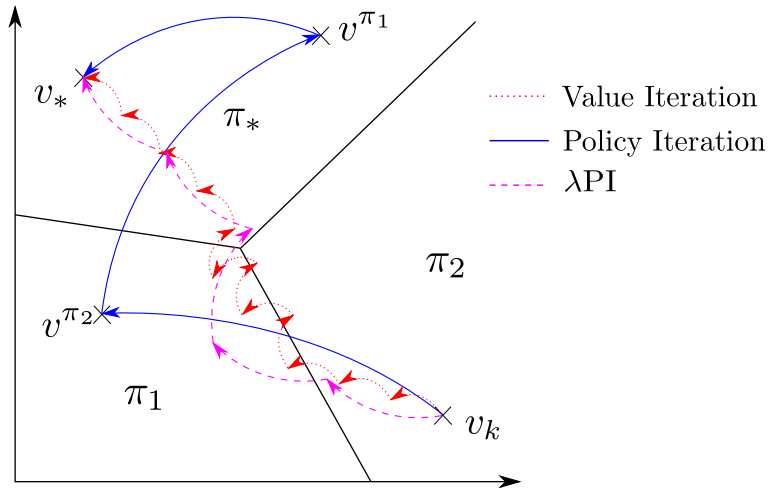


Figure 1: **Visualizing  $\lambda$  Policy Iteration on the greedy partition sketch:** Following (Bertsekas and Tsitsiklis, 1996, page 226), one can decompose the value space as a collection of polyhedra, such that each polyhedron corresponds to a region where one policy is greedy. This is called the *greedy partition*. In the above example, there are only 3 policies,  $\pi_1$ ,  $\pi_2$  and  $\pi_*$ .  $v_k$  is the initial value.  $\text{greedy}(v_k) = \pi_2$ ,  $\text{greedy}(v^{\pi_2}) = \pi_1$ , and  $\text{greedy}(v^{\pi_1}) = \pi_*$ . Therefore (1-)Policy Iteration generates the sequence  $((\pi_2, v^{\pi_2}), (\pi_1, v^{\pi_1}), (\pi_*, v^{\pi_*}))$ . Value iteration (or 0 Policy Iteration) starts by slowly updating  $v_k$  towards  $v^{\pi_2}$  until it crosses the boundary  $\pi_1/\pi_2$ , after which it tracks alternatively  $v^{\pi_1}$  and  $v^{\pi_2}$ , until it reaches the  $\pi_*$  part. In other words, Value Iteration makes small steps.  $\lambda$  Policy Iteration is intermediate between the two: it makes steps of which the steps length is related to  $\lambda$ . On the above sketch,  $\lambda$  Policy Iteration gets to the greedy partition of the optimal policy  $\pi_*$  in 3 steps, as did Policy Iteration.

- $M_k$  is a contraction mapping of modulus  $\lambda\gamma$  for the max-norm ;
- The next iterate  $v_{k+1}$  of  $\lambda$  Policy Iteration is the (unique) fixed point of  $M_k$ .

The left representation of  $\lambda$  Policy Iteration is obtained by “unrolling” Equation 8 an infinite number of times, while the right one is obtained by using Equation 9 and solving the linear system  $v_{k+1} = M_k v_{k+1}$ .

Informally, the parameter  $\lambda$  (or  $n$  in the case of Modified Policy Iteration) can be seen as adjusting the size of the step for tracking the target  $v^{\pi_{k+1}}$  (see Figure 1): the bigger the value, the longer the step. Formally,  $\lambda$  Policy Iteration (consider the above left hand side) consists in doing a  $\lambda$ -geometric average of the different number of applications of the Bellman operator  $(\mathcal{T}^{\pi_{k+1}})^j$  to  $v_k$ . The right hand side is here interesting because it clearly shows that  $\lambda$  Policy Iteration generalizes Value Iteration (when  $\lambda = 0$ ) and Policy Iteration (when  $\lambda = 1$ ). The operator  $M_k$  gives some insight on how one may concretely implement one iteration of  $\lambda$  Policy Iteration: it can for instance be done through a Value Iteration-like

algorithm which applies  $M_k$  iteratively. Also, the fact that its contraction factor  $\lambda\gamma$  can be tuned is of particular importance because finding the corresponding fixed point can be much easier than that of  $\mathcal{T}^{\pi_{k+1}}$ , which is only  $\gamma$ -contracting.

In order to describe the  $\lambda$  Policy Iteration algorithm, it is useful to introduce a (last) new operator. For any value  $v$  and any policy  $\pi$ , define<sup>9</sup>:

$$\mathcal{T}_\lambda^\pi v := v + (I - \lambda\gamma P^\pi)^{-1}(\mathcal{T}^\pi v - v) \quad (10)$$

$$= (I - \lambda\gamma P^\pi)^{-1}(\mathcal{T}^\pi v - \lambda\gamma P^\pi v) \quad (11)$$

$$= (I - \lambda\gamma P^\pi)^{-1}(r^\pi + (1 - \lambda)\gamma P^\pi v) \quad (11)$$

$$= (I - \lambda\gamma P^\pi)^{-1}(\lambda r^\pi + (1 - \lambda)\mathcal{T}^\pi v) \quad (12)$$

$\lambda$  Policy Iteration is formally described in Algorithm 3. Once again, our description

---

### Algorithm 3 $\lambda$ Policy Iteration

---

**Input:** An MDP,  $\lambda \in (0, 1)$ , an initial value  $v_0$

**Output:** An (approximately) optimal policy

$k \leftarrow 0$

**repeat**

$\pi_{k+1} \leftarrow \text{greedy}(v_k)$       // Update the policy

$v_{k+1} \leftarrow \mathcal{T}_\lambda^{\pi_{k+1}} v_k + \epsilon_{k+1}$       // Update the estimate of the value of policy  $\pi_{k+1}$

$k \leftarrow k + 1$

**until** some convergence criterion

**Return**  $\text{greedy}(v_k)$

---

includes a potential error term each time the value is updated. Even with this error term, it is straightforward to see that the algorithm reduces to Value Iteration (Algorithm 1) when  $\lambda = 0$  and to Policy Iteration<sup>10</sup> (Algorithm 2) when  $\lambda = 1$ .

**Relation with Reinforcement Learning** The definition of the operator  $\mathcal{T}_\lambda^\pi$  given by Equation 11 is the form we have used for introducing of  $\lambda$  Policy Iteration as an intermediate algorithm between Value Iteration and Policy Iteration. The equivalent form given by Equation 10 can be used to make a connection with the TD( $\lambda$ ) algorithms<sup>11</sup> (Sutton and Barto, 1998). Indeed, through Equation 10, the evaluation phase of  $\lambda$  Policy Iteration can be seen as an incremental additive procedure:

$$v_{k+1} \leftarrow v_k + \Delta_k$$

where

$$\Delta_k := (I - \lambda\gamma P^{\pi_{k+1}})^{-1}(\mathcal{T}^{\pi_{k+1}} v_k - v_k)$$

---

9. The four formulations are equivalent up to some little linear algebra manipulations.

10. Policy Iteration starts with an initial policy while  $\lambda$  Policy Iteration starts with some initial value. To be precise, 1 Policy Iteration starting with  $v_0$  is equivalent to Policy Iteration starting with the greedy policy with respect to  $v_0$ .

11. TD stands for Temporal Difference. As we have mentioned in Footnote 7,  $\lambda$  Policy Iteration was originally also called ‘‘Temporal Difference Based Policy Iteration’’ and the presentation of Bertsekas and Ioffe (1996) starts from the formulation of Equation 10 (which is close to TD( $\lambda$ )), and afterwards makes the connection with Value Iteration and Policy Iteration.

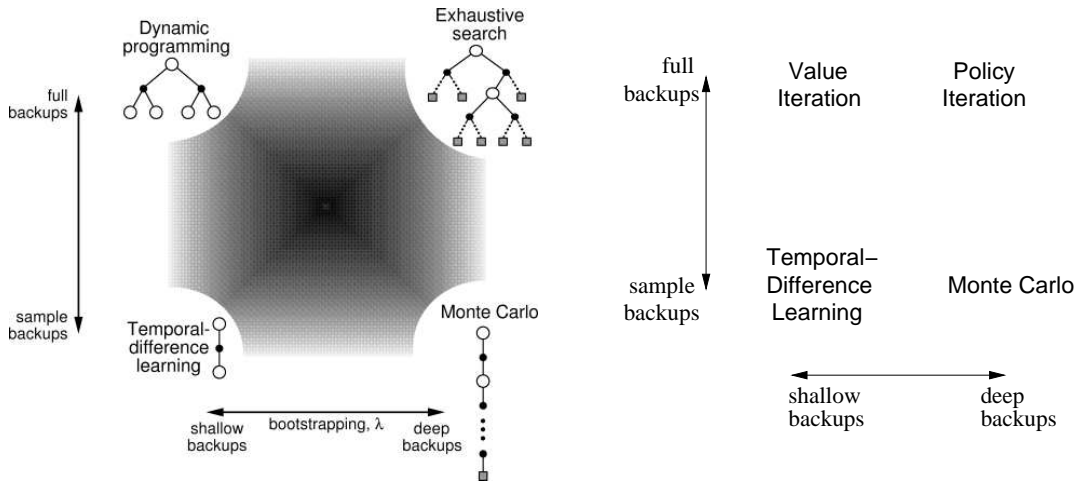


Figure 2:  **$\lambda$  Policy Iteration, a fundamental algorithm for Reinforcement Learning:** The left drawing, taken from chapter 10.1 the book of Sutton and Barto (1998), represents “two of the most important dimensions” of Reinforcement Learning methods. The vertical axis corresponds to whether one does full backup (exact computation of the expectations) or stochastic approximation (estimation through samples). The horizontal axis corresponds to the depth of the backups, and is (among other things) controlled by the parameter  $\lambda$ . On the right, is an original picture of  $\lambda$  Policy Iteration, along the same dimensions. It is interesting to notice that Sutton and Barto (1998) comment their drawing as follows: “At three of the four corners of the space are the three primary methods for estimating values: DP, TD, and Monte Carlo”, i.e. they do not recognize the fourth corner as one of the Reinforcement Learning *primary methods*. The natural representation of  $\lambda$  Policy Iteration actually suggests a modification of the sketch which is particularly meaningful: Policy Iteration, which consists in computing the value of the current policy, is the *deepest backup method*, and can be considered as the batch version of Monte Carlo.

is zero if and only if the value  $v_k$  is equal to the optimal value  $v_*$ . It can be shown (see Bertsekas and Ioffe (1996) for a proof or simply look at the equivalence between Equations 2 and 3 for the intuition) that the vector  $\delta_k$  has components given by:

$$\Delta_k(i) = \lim_{N \rightarrow \infty} E_{\pi_{k+1}} \left[ \sum_{j=0}^{N-1} (\lambda\gamma)^j \delta_k(i_j, i_{j+1}) \mid i_0 = i \right] \quad (13)$$

with

$$\delta_k(i, j) := r(i, \pi_{k+1}(i), j) + \gamma v(j) - v(i)$$

being the temporal difference associated to transition  $i \rightarrow j$ , as defined by Sutton and Barto (1998). When one uses a stochastic approximation of  $\lambda$  Policy Iteration, that is when the expectation  $E_{\pi_{t+1}}$  is approximated by sampling,  $\lambda$  Policy Iteration reduces to the algorithm TD( $\lambda$ ) which is described in chapter 7 of Sutton and Barto (1998). In particular, it is

well known that when  $\lambda = 1$ , the terms in the above sum collapse and become the exact discounted return:

$$\begin{aligned} \sum_{j=0}^{N-1} \gamma^j \delta_k(i_j, i_{j+1}) &= \sum_{j=0}^{N-1} \gamma^j [r(i_j, \pi_{k+1}(i_j), i_{j+1}) + \gamma v(i_{j+1}) - v(i_j)] \\ &= \sum_{j=0}^{N-1} \gamma^j r(i_j, \pi_{k+1}(i_j), i_{j+1}) + \gamma^N v(i_N) \end{aligned}$$

and the stochastic approximation matches the Monte-Carlo method. Also, Bertsekas and Ioffe (1996) show that Approximate TD( $\lambda$ ) with a linear feature architecture, as described in chapter 8.2 of Sutton and Barto (1998), corresponds to a natural Approximate version of  $\lambda$  Policy Iteration where the value is updated by least square fitting using a gradient-type iteration after each sample. Last but not least, the reader might notice that the “unified view” of Reinforcement Learning algorithms which is depicted in chapter 10.1 of Sutton and Barto (1998), and which is reproduced in Figure 2, is in fact a picture of  $\lambda$  Policy Iteration.

### 3.2 Known Properties of $\lambda$ Policy Iteration

To our knowledge, little has been done concerning the analysis of  $\lambda$  Policy Iteration: the only results available concern the Exact case (when  $\epsilon_k = 0$ ). Define the following factor

$$\beta = \frac{(1 - \lambda)\gamma}{1 - \lambda\gamma}. \quad (14)$$

We have  $0 \leq \beta \leq \gamma < 1$ . If  $\lambda = 0$  (Value Iteration) then  $\beta = \gamma$ , and if  $\lambda = 1$  (Policy Iteration) then  $\beta = 0$ . In the original article introducing  $\lambda$  Policy Iteration, Bertsekas and Ioffe (1996) show the convergence and provide an asymptotic rate of convergence:

**Proposition 11 (Convergence of Exact  $\lambda$ PI (Bertsekas and Ioffe, 1996))**

*If the discount factor  $\gamma < 1$ , then  $v_k$  converges to  $v_*$ . Furthermore, after some index  $k_*$ , the rate of convergence is linear in  $\beta$  as defined in Equation 14, that is*

$$\forall k \geq k_*, \quad \|v_{k+1} - v_*\| \leq \beta \|v_k - v_*\|.$$

By making  $\lambda$  close to 1,  $\beta$  can be arbitrarily close to 0 so the above rate of convergence might look overly impressive. This needs to be put into perspective: the index  $k_*$  is the index after which the policy  $\pi_k$  does not change anymore (and is equal to the optimal policy  $\pi_*$ ). As we said when we introduced the algorithm,  $\lambda$  controls the speed at which one wants  $v_k$  to “track the target”  $v^{\pi_{k+1}}$ ; when  $\lambda = 1$ , this is done in one step (and if  $\pi_{k+1} = \pi_*$  then  $v_{k+1} = v_*$ ). However, the bigger the value of  $\lambda$ , the less the operator  $M_k$  is contracting (recall that its contraction factor is  $\lambda\gamma$ ), and the harder the computation of its fixed point  $v_{k+1}$ .

## 4. Overview of our Results and Main Proof Ideas

Now that we have described the algorithms and some of their known properties, motivating the remaining of this paper is straightforward.  $\lambda$  Policy Iteration is conceptually nice since



it generalizes the two most-well known algorithms for solving discounted infinite-horizon Markov Decision Processes. The natural question that arises is whether one can generalize the results we have described so far to  $\lambda$  Policy Iteration (uniformly for all  $\lambda$ ). The answer is yes:

- We will derive a componentwise analysis of Exact and Approximate  $\lambda$  Policy Iteration.
- We will characterize the *non-asymptotic* convergence rate of Exact  $\lambda$  Policy Iteration (Proposition 11 only showed the *asymptotic* linear convergence) and will generalize the stopping criterion described for Value Iteration (Proposition 1).
- We will give bounds of the asymptotic error of Approximate  $\lambda$  Policy Iteration with respect to the *asymptotic* approximation error, Bellman residual, and Policy Bellman residual. This will generalize Lemmas 2 and 6, and Propositions 3, 4, 7 and 8. Our analysis actually implies that doing Approximate  $\lambda$  Policy Iteration is sound (when the approximation error tends to 0, the algorithm recovers the optimal solution).
- We will provide specific (better) bounds for the case when the value or the policy converges, which will thus generalize Corollaries 5 and 9.
- Interestingly, we will provide all our results using the span seminorms we have introduced at the beginning of the paper, and using the relations between this span semi-norms and the standard  $L_p$  norms (Equation 1), it can be seen that our results are in this respect slightly stronger than all the previously described results.

Conceptually, we will provide a unified vision (unified proofs, unified results) for all the mentioned algorithms.

#### 4.1 On the need for a new proof technique

In the literature, lines of analysis are different for Value Iteration and Policy Iteration. Analyses of Value Iteration are based on the fact that it computes the fixed point of the Bellman operator which is a  $\gamma$ -contraction mapping in max-norm (see e.g. Bertsekas and Tsitsiklis (1996)). Unfortunately, it can be shown that the operator by which Policy Iteration updates the value from one iteration to the next is in general not a contraction in max-norm. In fact, this observation can be drawn for  $\lambda$  Policy Iteration as soon as it does not reduce to Value Iteration:

**Proposition 12** *If  $\lambda > 0$ , there exists no norm for which the operator by which  $\lambda$  Policy Iteration updates the value from one iteration to the next is a contraction.*

**Proof** To see this, consider the following deterministic MDP with two states  $\{1, 2\}$  and two actions  $\{change, stay\}$ :  $r_1 = 0$ ,  $r_2 = 1$ ,  $P_{change}(s_2|s_1) = P_{change}(s_1|s_2) = P_{stay}(s_1|s_1) = P_{stay}(s_2|s_2) = 1$ . Consider the following two value functions  $v = (\epsilon, 0)$  and  $v' = (0, \epsilon)$  with  $\epsilon > 0$ . Their corresponding greedy policies are  $\pi = (stay, change)$  and  $\pi' = (change, stay)$ .

Then, we can compute the next iterates of  $v$  and  $v'$  (using Equation 11):

$$\begin{aligned} r^\pi + (1 - \lambda\gamma)P^\pi v &= ((1 - \lambda)\gamma\epsilon, 1 + (1 - \lambda)\gamma\epsilon) \\ \mathcal{T}_\lambda^\pi v &= \left( \frac{(1 - \lambda)\gamma\epsilon}{1 - \lambda\gamma}, 1 + \frac{(1 - \lambda)\gamma\epsilon}{1 - \lambda\gamma} \right) \\ r^{\pi'} + (1 - \lambda\gamma)P^{\pi'} v' &= ((1 - \lambda)\gamma\epsilon, 1 + (1 - \lambda)\gamma\epsilon) \\ \mathcal{T}_\lambda^{\pi'} v' &= \left( \frac{1 + (1 - \lambda)\gamma\epsilon}{1 - \lambda\gamma} - 1, \frac{1 + (1 - \lambda)\gamma\epsilon}{1 - \lambda\gamma} \right) \end{aligned}$$

Then

$$\mathcal{T}_\lambda^{\pi'} v' - \mathcal{T}_\lambda^\pi v = \left( \frac{1}{1 - \lambda\gamma} - 1, \frac{1}{1 - \lambda\gamma} - 1 \right)$$

while

$$v' - v = (-\epsilon, \epsilon).$$

As  $\epsilon$  can be arbitrarily small, the norm of  $\mathcal{T}_\lambda^\pi v - \mathcal{T}_\lambda^{\pi'} v'$  can be arbitrarily larger than norm of  $v - v'$  when  $\lambda > 0$ .  $\blacksquare$

Analyses of Policy Iteration usually rely on the fact that the sequence of values generated by Exact Policy Iteration is non-decreasing (see Bertsekas and Tsitsiklis (1996); Munos (2003)). Unfortunately, it can easily be seen that as soon as  $\lambda$  is smaller than 1, the value functions may decrease (it suffices to take a very high initial value). For non trivial values of  $\lambda$ ,  $\lambda$  Policy Iteration is neither contracting nor non-decreasing, so we need a new proof technique.

#### 4.2 An Overview on the Componentwise Analysis of $\lambda$ Policy Iteration

The rest of this section provides an overview of our analysis. We show how to compute an upper bound of the loss for  $\lambda$  Policy Iteration in the general (possibly approximate) case. It will be the basis for the derivation of componentwise bounds for Exact  $\lambda$  Policy Iteration (Section 5) and Approximate  $\lambda$  Policy Iteration (Section 6). Consider  $\lambda$  Policy Iteration as described in Algorithm 3, and the sequences of value-policy-error triplets  $(v_k, \pi_k, \epsilon_k)$  it generates. Most of our results come from a series of relations involving objects we now define:

- the **loss** of using policy  $\pi_k$  instead of the optimal policy:

$$l_k := v_* - v^{\pi_k}$$

- the **value** of the  $k^{\text{th}}$  iterate b.a. (before approximation):

$$w_k := v_k - \epsilon_k$$

- the **distance** between the optimal value and the  $k^{\text{th}}$  value b.a. :

$$d_k := v_* - w_k$$

- the **shift** between the  $k^{\text{th}}$  value b.a. and the value of the  $k^{\text{th}}$  policy:

$$s_k := w_k - v^{\pi_k}$$

- the **Bellman residual** of the  $k^{\text{th}}$  value:

$$b_k := \mathcal{T}_{k+1}v_k - v_k = \mathcal{T}v_k - v_k$$

To lighten the notations, we will now on write:  $P_k := P^{\pi_k}$ ,  $\mathcal{T}_k := \mathcal{T}^{\pi_k}$ ,  $P_* := P^{\pi_*}$ . We will refer to the factor  $\beta$  as introduced by Bertsekas and Ioffe (Equation 14 page 15). Also, the following stochastic matrix will play a recurrent role in our analysis<sup>12</sup> :

$$A_k := (1 - \lambda\gamma)(I - \lambda\gamma P_k)^{-1}P_k. \quad (15)$$

We will use the notation  $\bar{x}$  for an upper bound of  $x$  and  $\underline{x}$  for a lower bound.

Our analysis relies on a series of Lemmas that we now state (for clarity, all the proofs are deferred to Appendix A).

**Lemma 13** *The shift is related to the Bellman residual as follows:*

$$s_k = \beta(I - \gamma P_k)^{-1}A_k(-b_{k-1}).$$

**Lemma 14** *The Bellman residual at iteration  $k+1$  cannot be much lower than the Bellman residual at iteration  $k$ :*

$$b_{k+1} \geq \beta A_{k+1}b_k + x_{k+1}$$

where  $x_k := (\gamma P_k - I)\epsilon_k$  only depends on the approximation error.

As a consequence, a lower bound of the Bellman residual is<sup>13</sup>:

$$b_k \geq \sum_{j=k_0+1}^k \beta^{k-j} (A_k A_{k-1} \dots A_{j+1}) x_j + \beta^{k-k_0} (A_k A_{k-1} \dots A_{k_0+1}) b_{k_0} := \underline{b}_k.$$

Using Lemma 13, the bound on the Bellman residual also provides an upper on the shift<sup>14</sup>:

$$s_k \leq \beta(I - \gamma P_k)^{-1}A_k(-\underline{b}_{k-1}) := \bar{s}_k$$

**Lemma 15** *The distance tends to reduce:*

$$d_{k+1} \leq \gamma P_* d_k + y_k$$

where  $y_k := \frac{\lambda\gamma}{1-\lambda\gamma}A_{k+1}(-\underline{b}_k) - \gamma P_* \epsilon_k$  depends on the lower bound of the Bellman residual and the approximation error.

12. The fact that this is indeed a stochastic matrix is explained at the beginning of the Appendices.

13. We use the property here that if some vectors satisfy the componentwise inequality  $x \leq y$ , and if  $P$  is a stochastic matrix, then the componentwise inequality  $Px \leq Py$  holds.

14. We use the fact that  $(1 - \gamma)(I - \gamma P_k)^{-1}$  is a stochastic matrix (cf Footnote 12) and Footnote 13.

Then, an upper bound of the distance is<sup>15</sup>:

$$d_k \leq \sum_{j=k_0}^{k-1} \gamma^{k-1-j} (P_*)^{k-1-j} y_j + \gamma^{k-k_0} (P_*)^{k-k_0} d_{k_0} = \bar{d}_k.$$

Eventually, as

$$l_k = d_k + s_k \leq \bar{d}_k + \bar{s}_k,$$

the upper bounds on the distance and the shift will enable us to derive the upper bound on the loss.

**Remark 16** *The above proof is a generalization of that of Munos (2003) for Approximate Policy Iteration. Note that when  $\lambda = 1$ , that is when both proofs coincide,  $\beta = 0$  and Lemmas 13 and 14 have the following particularly trivial form:  $s_k = 0$  and  $b_{k+1} \geq x_{k+1}$ .*

The next two sections contain our main results, which take the form of bounds performance bounds when using  $\lambda$  Policy Iteration. Section 5 gathers the results concerning Exact  $\lambda$  Policy Iteration, while Section 6 presents those concerning Approximate  $\lambda$  Policy Iteration.

## 5. Performance bounds of Exact $\lambda$ Policy Iteration

Consider Exact  $\lambda$  Policy Iteration for which we have  $\epsilon_k = 0$  for all  $k$ . By exploiting the recursive relations we have described in the previous section (this process is detailed in Appendix B), we can derive the following componentwise bounds for the loss:

### Lemma 17 (Componentwise Rate of Convergence of Exact $\lambda$ PI)

*The following matrices*

$$\begin{aligned} E_{kk_0} &:= (1 - \gamma)(P_*)^{k-k_0}(I - \gamma P_*)^{-1} \\ E'_{kk_0} &:= \left( \frac{1 - \gamma}{\gamma^{k-k_0}} \right) \left( \frac{\lambda\gamma}{1 - \lambda\gamma} \sum_{j=k_0}^{k-1} \gamma^{k-1-j} \beta^{j-k_0} (P_*)^{k-1-j} A_{j+1} A_j \dots A_{k_0+1} \right. \\ &\quad \left. + \beta^{k-k_0} (I - \gamma P_k)^{-1} A_k A_{k-1} \dots A_{k_0+1} \right) \\ F_{kk_0} &:= (1 - \gamma) P_*^{k-k_0} + \gamma E'_{kk_0} P_* \end{aligned}$$

*are stochastic and the performance of the policies generated by  $\lambda$  Policy Iteration satisfies*

$$\begin{aligned} v_* - v^{\pi_k} &\leq \frac{\gamma^{k-k_0}}{1 - \gamma} [F_{kk_0} - E'_{kk_0}] (v_* - v_{k_0}) \\ v_* - v^{\pi_k} &\leq \frac{\gamma^{k-k_0}}{1 - \gamma} [E_{kk_0} - E'_{kk_0}] (\mathcal{T}v_{k_0} - v_{k_0}). \\ v_* - v^{\pi_k} &\leq \gamma^{k-k_0} \left[ (P_*)^{k-k_0} \left( (v_* - v_{k_0}) - \min_s [v_*(s) - v_{k_0}(s)] e \right) + \|v_*(s) - v^{\pi_{k_0+1}}\|_\infty e \right]. \end{aligned}$$

---

15. See Footnote 13.

In order to derive (more interpretable) span seminorms bounds from the above componentwise bound, we will rely on the following Lemma, which for clarity of exposition is proved in Appendix E.

**Lemma 18** *If for some non-negative vectors  $x$  and  $y$ , some constant  $K \geq 0$ , and some stochastic matrices  $X$  and  $X'$  we have*

$$x \leq K(X - X')y,$$

*Then*

$$\|x\|_\infty \leq K \text{span}_\infty [y].$$

With this, the componentwise bounds of Lemma 17 become:

**Proposition 19 (Non-asymptotic bounds for Exact  $\lambda$  Policy Iteration)**

*For any  $k > k_0$ ,*

$$\begin{aligned} \|v_* - v^{\pi_k}\|_\infty &\leq \frac{\gamma^{k-k_0}}{1-\gamma} \text{span}_\infty [v_* - v_{k_0}] \\ \|v_* - v^{\pi_k}\|_\infty &\leq \frac{\gamma^{k-k_0}}{1-\gamma} \text{span}_\infty [\mathcal{T}v_{k_0} - v_{k_0}] \\ \|v_* - v^{\pi_k}\|_\infty &\leq \gamma^{k-k_0} (\text{span}_\infty [v_* - v_{k_0}] + \|v_*(s) - v^{\pi_{k_0+1}}\|_\infty). \end{aligned} \quad (16)$$

This *non-asymptotic* bound supplements the *asymptotic* bound of Proposition 11 from Bertsekas and Ioffe (1996). The first bound is expressed in terms of the distance between the value function and the optimal value function at some iteration  $k_0$ . The second inequality can be used as a stopping criterion. Indeed, taking  $k = k_0 + 1$  it implies the following stopping condition, which generalizes that of Proposition 1 about Value Iteration:

**Proposition 20 (Stopping condition of Exact  $\lambda$  Policy Iteration)**

*If at some iteration  $k_0$ , the value  $v_{k_0}$  satisfies:*

$$\text{span}_\infty [\mathcal{T}v_{k_0} - v_{k_0}] \leq \frac{1-\gamma}{\gamma} \epsilon$$

*then the greedy policy  $\pi_{k_0+1}$  with respect to  $v_{k_0}$  is  $\epsilon$ -optimal:  $\|v_* - v^{\pi_{k_0+1}}\|_\infty < \epsilon$ .*

The last inequality relies on the distance between the value function and the optimal value function and the value difference between the optimal policy and the first greedy policy; compared to the others, it has the advantage of not containing a  $\frac{1}{1-\gamma}$  factor. To our knowledge, this bound is even new for the specific cases of Value Iteration and Policy Iteration.

## 6. Performance bounds for Approximate $\lambda$ Policy Iteration

We now turn to the (slightly more involved) results on Approximate  $\lambda$  Policy Iteration. We will provide componentwise bounds of the loss  $v_* - v^{\pi_k} \geq 0$  of using policy  $\pi_k$  instead of using the optimal policy, with respect to the approximation error  $\epsilon_k$ , the Policy Bellman

residual  $\mathcal{T}_k v_k - v_k$  and the Bellman residual  $\mathcal{T} v_k - v_k = \mathcal{T}_{k+1} v_k - v_k$ . Recall the subtle difference between these two Bellman residuals: the Policy Bellman residual says how much  $v_k$  differs from the value of  $\pi_k$  while the Bellman residual says how much  $v_k$  differs from the value of the policies  $\pi_{k+1}$  and  $\pi_*$ .

The core of our analysis is based on the following Lemma:

**Lemma 21 (Componentwise Performance bounds for App.  $\lambda$  Policy Iteration)**

The following matrices

$$\begin{aligned}
 B_{jk} &:= \frac{1-\gamma}{\gamma^{k-j}} \left[ \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} (P_*)^{k-1-i} A_{i+1} A_i \dots A_{j+1} \right. \\
 &\qquad \qquad \qquad \left. + \beta^{k-j} (I - \gamma P_k)^{-1} A_k A_{k-1} \dots A_{j+1} \right] \\
 B'_{jk} &:= \gamma B_{jk} P_j + (1-\gamma)(P_*)^{k-j} \\
 C_k &:= (1-\gamma)^2 (I - \gamma P_*)^{-1} (P_*(I - \gamma P_k)^{-1}) \\
 C'_k &:= (1-\gamma)^2 (I - \gamma P_*)^{-1} (P_{k+1}(I - \gamma P_{k+1})^{-1}) \\
 D &:= (1-\gamma) P_* (I - \gamma P_*)^{-1} \\
 D'_k &:= (1-\gamma) P_k (I - \gamma P_k)^{-1}
 \end{aligned}$$

are stochastic and

$$\begin{aligned}
 \forall k_0, \quad \limsup_{k \rightarrow \infty} v_* - v^{\pi_k} &\leq \frac{1}{1-\gamma} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} [B_{jk} - B'_{jk}] \epsilon_j \\
 \limsup_{k \rightarrow \infty} v_* - v^{\pi_k} &\leq \frac{\gamma}{(1-\gamma)^2} \limsup_{k \rightarrow \infty} [C_k - C'_k] (\mathcal{T}_k v_k - v_k) \\
 \forall k, \quad v_* - v^{\pi_k} &\leq \frac{\gamma}{1-\gamma} [D - D'_k] (\mathcal{T} v_{k-1} - v_{k-1}).
 \end{aligned}$$

The first inequality involving the errors ( $\epsilon_k$ ) is based on Lemmas 13-15 (presented in Section 4) and is detailed in Appendix C. The two other inequalities (the asymptotic performance of Approximate  $\lambda$  Policy Iteration with respect to the Bellman residuals) is somewhat simpler and is proved independently in Appendix D.

**Remark 22 (Relation with the previous bounds of Munos (2007, 2003))** We can look at the relation between our bound for general  $\lambda$  and the bounds derived separately by Munos for Approximate Value Iteration (Lemma 2) and Approximate Policy Iteration (Lemma 6):

- Let us first consider the case where  $\lambda = 0$ . Then  $\beta = \gamma$ ,  $A_k = P_k$  and

$$B_{jk} = (1-\gamma)(I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_{j+1}.$$

Then our bound implies that  $\limsup_{k \rightarrow \infty} v_* - v^{\pi_k}$  is upper bounded by:

$$\begin{aligned} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} \left[ (I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_{j+1} \right. \\ \left. - \left( \gamma (I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_j + (P_*)^{k-j} \right) \right] \epsilon_j. \end{aligned} \quad (17)$$

The bound derived by Munos for Approximate Value Iteration (Lemma 2 page 6) is

$$\begin{aligned} \limsup_{k \rightarrow \infty} (I - \gamma P_k)^{-1} \sum_{j=0}^{k-1} \gamma^{k-j} \left[ P_k P_{k-1} \dots P_{j+1} - (P_*)^{k-j} \right] \epsilon_j \\ = \limsup_{k \rightarrow \infty} \sum_{j=0}^{k-1} \gamma^{k-j} \left[ (I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_{j+1} - (I - \gamma P_k)^{-1} (P_*)^{k-j} \right] \epsilon_j \\ = \limsup_{k \rightarrow \infty} \sum_{j=0}^{k-1} \gamma^{k-j} \left[ (I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_{j+1} \right. \\ \left. - \left( (I - \gamma P_k)^{-1} \gamma P_k (P_*)^{k-j} + (P_*)^{k-j} \right) \right] \epsilon_j. \end{aligned} \quad (18)$$

The above bounds are very close to each other: we can go from Equation 17 to Equation 18 by replacing  $P_{k-1} \dots P_j$  by  $(P_*)^{k-j}$ .

- Now, when  $\lambda = 1$ ,  $\beta = 0$ ,  $A_k = (1 - \gamma)(I - \gamma P_k)^{-1} P_k$  and

$$B_{jk} = (1 - \gamma)(P_*)^{k-1-j} P_{j+1} (I - \gamma P_{j+1})^{-1}.$$

Our bound is

$$\limsup_{k \rightarrow \infty} v_* - v^{\pi_k} \leq \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-1-j} (P_*)^{k-1-j} u_j$$

with

$$u_j := [\gamma P_{j+1} (I - \gamma P_{j+1})^{-1} (I - \gamma P_j) - \gamma P_*] \epsilon_j.$$

By definition of the supremum limit, for all  $\epsilon > 0$ , there exists an index  $k_1$  such that for all  $j \geq k_1$ ,

$$u_j \leq \limsup_{l \rightarrow \infty} u_l + \epsilon.$$

Then:

$$\begin{aligned} \limsup_{k \rightarrow \infty} \sum_{j=k_1}^{k-1} \gamma^{k-1-j} (P_*)^{k-1-j} u_j &\leq \limsup_{k \rightarrow \infty} \sum_{j=k_1}^{k-1} \gamma^{k-1-j} (P_*)^{k-1-j} \left( \limsup_{l \rightarrow \infty} u_l + \epsilon \right) \\ &= (I - \gamma P_*)^{-1} \left( \limsup_{l \rightarrow \infty} u_l + \epsilon \right). \end{aligned}$$

As this is true for all  $\epsilon > 0$ , we eventually find the bound of Munos for Approximate Policy Iteration (Lemma 6 page 9).

Thus, up to some little details, our componentwise analysis unifies those of Munos. It is not a surprise that we fall back on the result of Munos for Approximate Policy Iteration because, as already mentioned in Remark 16, the proof we developed in Section 4 and Appendix A is a generalization of his. If we don't exactly recover the componentwise analysis of Munos for Approximate Value Iteration, this is not really fundamental as most of the results we are going to deduce won't be affected.

The componentwise bounds on the performance of Approximate  $\lambda$  Policy Iteration can be translated into span seminorm bounds, using the following Lemma (proved in Appendix E):

**Lemma 23** *If  $x_k$  and  $y_k$  are sequences of vectors and  $X_{jk}$ ,  $X'_{jk}$  sequences of stochastic matrices satisfying*

$$\forall k_0, \quad \limsup_{k \rightarrow \infty} |x_k| \leq K \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \xi_{k-j} (X_{kj} - X'_{kj}) y_j,$$

where  $(\xi_i)_{i \geq 1}$  is a sequence of non-negative weights satisfying:

$$\sum_{i=1}^{\infty} \xi_i = K' < \infty,$$

then, for all distribution  $\mu$ ,

$$\mu_{kj} := \frac{1}{2} (X_{kj} + X'_{kj})^T \mu$$

are distributions and

$$\limsup_{k \rightarrow \infty} \|x_k\|_{p, \mu} \leq KK' \lim_{k_0 \rightarrow \infty} \left[ \sup_{k \geq j \geq k_0} \text{span}_{p, \mu_{kj}} [y_j] \right].$$

Thus, using this Lemma and the fact that  $\sum_{i=1}^{\infty} \gamma^i = \frac{\gamma}{1-\gamma}$ , Lemma 21 can be truned into the following Proposition that unifies and generalizes Proposition 3 (page 6) on Approximate Value Iteration and Proposition 7 (page 9) on Approximate Policy Iteration.

**Proposition 24 (Span Seminorm Asymptotic Performance of  $\lambda$ API (1/2))**

*With the notations of Lemma 21, for all  $p$ ,  $k \geq j \geq 0$  and all distribution  $\mu$ ,*

$$\begin{aligned} \mu_{kj} &:= \frac{1}{2} (B_{jk} + B'_{jk})^T \mu, \\ \mu'_{kj} &:= \frac{1}{2} (C_k + C'_k)^T \mu, \text{ and} \\ \mu''_{kj} &:= \frac{1}{2} (D + D'_k)^T \mu \end{aligned}$$

are distributions and the performance of the policies generated by  $\lambda$  Policy Iteration satisfy:

$$\begin{aligned} \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{p, \mu} &\leq \frac{\gamma}{(1-\gamma)^2} \lim_{k_0 \rightarrow \infty} \left[ \sup_{k \geq j \geq k_0} \text{span}_{p, \mu_{kj}} [\epsilon_j] \right], \\ \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{p, \mu} &\leq \frac{\gamma}{(1-\gamma)^2} \limsup_{k \rightarrow \infty} \text{span}_{p, \mu'_{kj}} [\mathcal{T}_k v_k - v_k], \\ \forall k, \quad \|v_* - v^{\pi_k}\|_{p, \mu} &\leq \frac{\gamma}{1-\gamma} \text{span}_{p, \mu''_{kj}} [\mathcal{T} v_{k-1} - v_{k-1}]. \end{aligned}$$



As already mentioned, a drawback of the above  $L_p$  bounds comes from the fact that the distributions involved on the right hand sides are unknown. To go round this issue, one may consider the **concentration coefficient** assumption introduced by Munos (2003, 2007) and already mentioned in Equation 6 page 7. For clarity, we recall its definition here. We assume there exists a distribution  $\nu$  and a real number  $C(\nu)$  such that

$$C(\nu) := \max_{i,j,a} \frac{p_{ij}(a)}{\nu(j)}.$$

Then, we have the following property:

**Lemma 25** *Let  $X$  be an average of products of stochastic matrices of the MDP. For any distribution  $\mu$ , and vector  $y$  and any  $p$ ,*

$$\text{span}_{p,\mu X} [y] \leq (C(\nu))^{1/p} \text{span}_{p,\nu} [y].$$

**Proof** It can be seen from the definition of the concentration coefficient  $C(\nu)$  that  $\mu^T X \leq C(\nu)\nu^T$ . Thus,

$$\begin{aligned} (\text{span}_{p,X^T\mu} [y])^p &= \min_a \left( \|y - ae\|_{p,X^T\mu} \right)^p \\ &= \min_a \mu^T X |y - ae|^p \\ &\leq C(\nu) \min_a \nu^T |y - ae|^p \\ &= C(\nu) \min_a \left( \|y - ae\|_{p,\nu} \right)^p \\ &= C(\nu) (\text{span}_{p,\nu} [y])^p. \quad \blacksquare \end{aligned}$$

Using this Lemma, and the fact that for any  $p$ ,  $\|x\|_\infty = \max_\mu \|x\|_{p,\mu}$ , the  $L_p$  bounds of Proposition 24 become

**Proposition 26 (Span Seminorm Asymptotic Performance of A $\lambda$ PI (2/2))**

*Let  $C(\nu)$  be the concentration coefficient defined in Equation 6. For all  $p$  and all  $k$ ,*

$$\begin{aligned} \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty &\leq \frac{\gamma}{(1-\gamma)^2} [C(\nu)]^{1/p} \limsup_{j \rightarrow \infty} \text{span}_{p,\nu} [\epsilon_j] \\ \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty &\leq \frac{\gamma}{(1-\gamma)^2} [C(\nu)]^{1/p} \limsup_{k \rightarrow \infty} \text{span}_{p,\nu} [\mathcal{T}_k v_k - v_k] \\ \forall k, \quad \|v_* - v^{\pi_k}\|_\infty &\leq \frac{\gamma}{1-\gamma} [C(\nu)]^{1/p} \text{span}_{p,\nu} [\mathcal{T} v_{k-1} - v_{k-1}]. \end{aligned}$$

This results generalizes and unifies those derived for Approximate Value Iteration (Proposition 4 page 7) and Approximate Policy Iteration (Proposition 8 page 9).

When comparing the specific bounds of Munos for Approximate Value Iteration (Propositions 3 and 4) and Approximate Policy Iteration (Propositions 7 and 8), we wrote that the latter had the nice property that the bounds only depend on *asymptotic* errors/residuals (while the former depends on all errors). Our bounds for  $\lambda$  Policy Iteration have this nice property too. Considering the relations between the span seminorms and the other standard norms (Equation 1 page 3), we see that our results are not only more general, but also slightly finer than those of Munos.

**When the policy or the value converges** The performance bounds with respect to the approximation error can be improved if we know or observe that the value or the policy converges. Note that the former condition implies the latter (while the opposite is not true: the policy may converge while the value still oscillates). Indeed, we have the following Corollary.

**Corollary 27 (Asymptotic Performance of A $\lambda$ PI in case of convergence)**

*If the value converges to some  $v$ , then the approximation error converges to some  $\epsilon$ , and the corresponding greedy policy  $\pi$  satisfies*

$$\|v_* - v^\pi\|_\infty \leq \frac{\gamma}{1-\gamma} [C(\nu)]^{1/p} \text{span}_{p,\nu}[\epsilon].$$

*If the policy converges to some  $\pi$ , then*

$$\|v_* - v^\pi\|_\infty \leq \frac{\gamma(1-\lambda\gamma)}{(1-\gamma)^2} [C(\nu)]^{1/p} \limsup_{j \rightarrow \infty} \text{span}_{p,\nu}[\epsilon_j].$$

The first bound is proved in Appendix C.2 while the second is proved in Appendix C.3. These bounds unify and extend those presented for Approximate Value Iteration (Corollary 5 page 7) and Approximate Policy Iteration (Corollary 9 page 10), in the similar situation where the policy or the value converges. It is interesting to notice that in the weaker situation where only the policy converges, the constant decreases from  $\frac{1}{(1-\gamma)^2}$  to  $\frac{1}{1-\gamma}$  when  $\lambda$  varies from 0 to 1; in other words, the closer to Policy Iteration, the better the bound in that situation.

## 7. Application of $\lambda$ Policy Iteration to the Game of Tetris

In the final part of this paper, we will consider (and describe for the sake of keeping this paper self-contained) exactly the same application (Tetris) and implementation as Bertsekas and Ioffe (1996). Our motivation for doing so is twofold:

- From a theoretical point of view, we will show how our analysis (made in the discounted case  $\gamma < 1$ ) can be adapted to an undiscounted problem (where  $\gamma = 1$ ) like Tetris.
- We obtained empirical results that are different (and much less intriguing) than those of the original study. This will give us the opportunity to describe what we think are the reasons for such a difference.

But before we do so, let us begin by describing the Tetris domain.

### 7.1 The Game of Tetris and its model as an MDP

Tetris is a popular video game created in 1985 by Alexey Pajitnov. The game is played on a  $10 \times 20$  grid where pieces of different shapes fall from the top (see Figure 3). The player has to choose where each piece is added: he can move it horizontally and rotate it. When a row is filled, it is removed and all cells above it move one row downwards. The goal is to

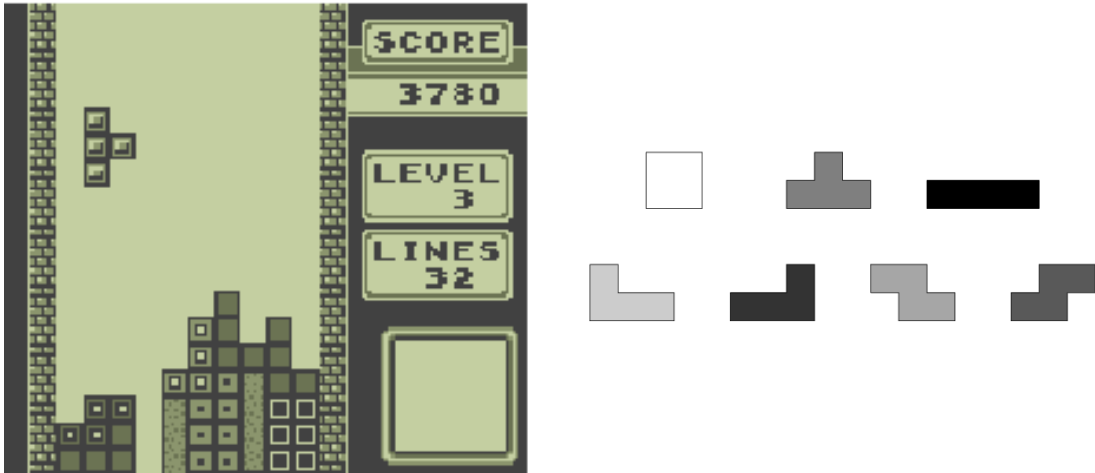


Figure 3: **Left:** a screenshot of a Tetris game. **Right:** the seven existing shapes.

remove as many lines as possible before the game is over, that is when there is not enough space remaining on the top of the pile to put the current new piece.

Instead of mimicking the original game (precisely described by Fahey (2003)), Bertsekas and Ioffe (1996) have focused on the main problem, that is choosing where to drop each coming piece. The corresponding MDP model is straightforward: The state consists of the wall configuration and the shape of the current falling piece. An action is the horizontal translation and the rotation which are applied to the piece before it is dropped on the wall. The reward is the number of lines which are removed after we have dropped the piece. As one considers the maximization of the score (the total number of lines removed during a game), the natural choice for the discount factor is  $\gamma = 1$ .

In a bit more details, the dynamics of Tetris is made of two components: the place where one drops the current piece and the choice of a new piece. As the latter component is uncontrollable (a new piece is chosen with uniform probability), the value functions needs not to be computed for all wall-piece pairs configurations but only for all wall configurations (see for instance Bertsekas and Ioffe (1996)). Also considering that the first component of the dynamics is deterministic, the optimal value function satisfies a reduced form of the Bellman Equation

$$\forall s \in S, v_*(s) = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \max_{a \in A(p)} r(s, p, a) + v_*(succ(s, p, a)) \quad (19)$$

where  $S$  is the set of wall configurations,  $\mathcal{P}$  is the set of pieces,  $A(p)$  is the set of translation-rotation pairs that can be applied to a piece  $p$ ,  $r(s, p, a)$  and  $succ(s, p, a)$  are respectively the number of lines removed and the (deterministic) next wall configuration if one puts a piece  $p$  on the wall  $s$  in translation-orientation  $a$ . The only function that satisfies the above Equation gives, for each wall configuration  $s$ , the average best score that can be achieved from  $s$ . If this function is known, a one step look-ahead strategy (i.e. a greedy policy) will perform optimally.

**Extension of the analysis for the undiscounted optimal control problem Tetris**

As just explained, the Tetris domain has a discount factor  $\gamma$  equal to 1, which makes it an *undiscounted MDP*. If this prevents us from applying directly most of the analysis we have made so far (since most of our bounds have a  $(1 - \gamma)$  term on the denominator), we are briefly going to show in what follows how to adapt the analysis so that we recover a significant error analysis.

In undiscounted infinite horizon control problems, it is generally assumed that there exists a  $N + 1^{\text{th}}$  termination state 0. Once the system reaches this state, it remains there for ever with no further reward, i.e. formally:

$$\forall a, p_{00}(a) = 1 \text{ and } r(0, a, 0) = 0.$$

In the case of Tetris, the termination state corresponds to “game over”, and the situation is particularly simple since Burgiel (1997) has shown that, whatever the strategy, some sequence of pieces (which will necessarily occur in finite time with probability 1) leads to game-over whatever the decisions taken<sup>16</sup>. This implies in particular that there exists an integer  $n_0 \leq N$  and a real number  $\alpha < 1$  such that for all initial distributions  $\mu$ , and actions  $a_0, a_1, \dots, a_{n_0-1}$ ,

$$P [i_{n_0} \neq 0 | i_0 \sim \mu, a_0, \dots, a_{n_0-1}] \leq \alpha. \quad (20)$$

We can define the MDP model for Tetris only on the  $N$  non-terminal states, that is on  $\{1, \dots, N\}$ . In this situation, for any policy  $\pi$ , the matrix  $P_\pi$  is in general a *substochastic* matrix (a matrix of which the components are non-negative and for which the max-norm is smaller than or equal to 1), and the above assumption means that for all set of  $n_0$  policies  $\pi_1, \pi_2, \dots, \pi_{n_0}$ ,

$$\|P_{\pi_1} P_{\pi_2} \cdots P_{\pi_{n_0}}\|_\infty \leq \alpha.$$

The componentwise analysis of  $\lambda$  Policy Iteration is here identical to what we have done before, except that we have  $\gamma = 1$  and  $\beta = 1$ <sup>17</sup>. The matrix  $A_k$  that appeared recurrently in our analysis has the following special form:

$$A_k := (1 - \lambda)(I - \lambda P_k)^{-1} P_k.$$

and is a substochastic matrix. The first bound of the componentwise analysis of  $\lambda$  Policy Iteration (Lemma 21 page 21) can be shown to be generalized as follows (see Appendix F for details):

**Lemma 28 (Componentwise Bounds in the Undiscounted Case)** *Assume that there exists  $n_0$  and  $\alpha$  such that Equation 20 holds. Write  $\eta := \frac{1 - \lambda^{n_0}}{1 - \lambda^{n_0} \alpha}$ . For all  $i$ , write*

$$\delta_i := \alpha \lfloor \frac{i}{n_0} \rfloor \left[ \left( \frac{1 - \lambda^{n_0}}{1 - \lambda} \right) \left( \frac{\lambda}{1 - \lambda^{n_0} \alpha} \right) \left( \frac{1 - \eta^i}{1 - \eta} \right) + \frac{n_0 \eta^i}{1 - \alpha} \right].$$

16. In the literature, a stationary policy that reaches the terminal state in finite time with probability 1 is said to be *proper*. The usual assumptions in undiscounted infinite horizon control problems are: 1) there exists at least one proper policy and 2) for every improper policy  $\pi$ , the corresponding value equals  $-\infty$  for at least one state. A simpler situation is when all stationary policies are proper. The situation of Tetris is even simpler: all (non necessarily stationary) policies are proper.

17. For simplicity in our discussion, we consider  $\lambda < 1$  (to avoid the special case  $\lambda = 1$  for which  $\beta$  is 0). The interested reader may however check that the results that we are going to state are continuous in the neighborhood of  $\lambda = 1$ .

Then the performance loss of the policies generated by  $\lambda$  Policy Iteration satisfy

$$\forall k_0, \limsup_{k \rightarrow \infty} v_* - v^{\pi_k} \leq \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \delta_{k-j} [G_{jk} - G'_{jk}] \epsilon_j \quad (21)$$

where for all  $j < k$ ,

$$G_{jk} := \frac{1}{\delta_{k-j}} \left[ \frac{\lambda}{1-\lambda} \sum_{i=j}^{k-1} (P_*)^{k-1-i} A_{i+1} A_i \dots A_{j+1} + (I - P_k)^{-1} A_k A_{k-1} \dots A_{j+1} \right], \text{ and}$$

$$G'_{jk} := \frac{1}{\delta_{k-j}} G_{jk} P_j$$

are substochastic matrices.

**Remark 29** By observing that  $\eta \in (0, 1)$ , and that for all  $x \in (0, 1)$ ,  $0 \leq \frac{1-x^{n_0}}{1-x} \leq n_0$ , it can be seen that the coefficients  $\delta_i$  are finite for all  $i$ . Furthermore, when  $n_0 = 1$  (which matches the discounted case with  $\alpha = \gamma$ ), one can observe that  $\delta_i = \frac{\gamma^i}{1-\gamma}$  and that one recovers the result of Lemma 21.

This Lemma can then be exploited to show that  $\lambda$  Policy Iteration enjoys an  $L_p$  norm guarantee. Indeed, an analogue of Proposition 24 (whose proof is detailed in Appendix F) is the following Proposition.

**Proposition 30 ( $L_p$  norm Bound for in the Undiscounted Case)**

Let  $C(\nu)$  be the concentration coefficient defined in Equation 6. Let the assumptions of Lemma 28 hold. Recall that we wrote  $\eta = \frac{1-\lambda^{n_0}}{1-\lambda^{n_0}\alpha}$ . For all distribution  $\mu$  on  $(1, \dots, N)$ ,  $\mu_{kj} := \frac{1}{2} (\tilde{G}_{jk} + \tilde{G}'_{jk})^T \mu$  is a non-negative vector and  $\tilde{\mu}_{kj} := \frac{\mu_{kj}}{\|\mu_{kj}\|_1}$  is a distribution on  $(1, \dots, N)$ . Then for all  $p$ , the loss of the policies generated by  $\lambda$  Policy Iteration satisfy

$$\limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{p, \mu} \leq K(\lambda, n_0) (C(\nu))^{1/p} \lim_{j \rightarrow \infty} \|\epsilon_j\|_{p, \tilde{\mu}_{kj}}$$

where

$$K(\lambda, n_0) := \lambda f(\lambda) \frac{f(1) - f(\eta)}{1 - \eta} + f(1) f(\eta) - f(1),$$

$$\forall x < 1, f(x) := \frac{(1 - x^{n_0})}{(1 - x)(1 - x^{n_0}\alpha)}, \text{ and by continuity } f(1) := \frac{n_0}{1 - \alpha}.$$

**Remark 31** There are three differences with respect to the results we have presented for the discounted case.

1. The fact that we defined the model (and thus the algorithm) only on the non-terminal states  $(1, \dots, N)$  implies that there is no error incurred in the terminal state 0. Note, however, that this is not a strong assumption since the value of the terminal state is necessarily 0.

2. The right hand side depends on the  $L_p$  norm, and not the span  $L_p$  seminorm. This is due to the fact that the matrices  $G_{jk}$  and  $G'_{jk}$  defined above are in general substochastic matrices (and not stochastic matrices).
3. Eventually, the constant  $K(\lambda, n_0)$  depends on  $\lambda$ . More precisely, it can be observed that:

$$\lim_{\lambda \rightarrow 0} K(\lambda, n_0) = \lim_{\lambda \rightarrow 1} K(\lambda, n_0) = \frac{n_0^2}{(1-\alpha)^2} - \frac{n_0}{1-\alpha}$$

and that this is the minimal value of  $\lambda \mapsto K(\lambda, n_0)$ . Though we took particular care in deriving this bound, we leave for future work the question whether one could provide a bound of  $\frac{n_0^2}{(1-\alpha)^2} - \frac{n_0}{1-\alpha}$  for any  $\lambda \in (0, 1)$ . When  $n_0 = 1$  (which matches the discounted case with  $\alpha = \gamma$ ),  $K(\lambda, 1)$  does not depend anymore on  $\lambda$  and we recover (without surprise) the bound of Proposition 24:

$$\forall \lambda, K(\lambda, 1) = \frac{\alpha}{(1-\alpha)^2}.$$

Now that we are reassured about the fact that applying  $\lambda$  Policy Iteration approximately to Tetris is principled, we turn to the precise description of its actual implementation.

## 7.2 An Instance of Approximate $\lambda$ Policy Iteration

For large scale problems, many Approximate Dynamic Programming algorithms are based on two complementary tricks:

- one uses samples to approximate the expectations such as that of Equation 13 ;
- one only looks for a linear approximation of the optimal value function:

$$v^r(s) = r(0) + \sum_{k=1}^K r(k) \Phi_k(s)$$

where  $r = (r(0) \dots r(K))$  is the parameter vector and  $\Phi_k(s)$  are some predefined feature functions on the state space. Thus, each value of  $r$  characterizes a value function  $v^r$  over the entire state space.

The instance of Approximate  $\lambda$  Policy Iteration of Bertsekas and Ioffe (1996) follows these ideas. More specifically, this algorithm is devoted to MDPs which have a termination state, that has 0 reward and is absorbing. For this algorithm to be run, one must further assume that all policies are proper, i.e. that all policies reach the termination state with probability one in finite time<sup>18</sup>.

Similarly to Exact  $\lambda$  Policy Iteration, this Approximate  $\lambda$  Policy Iteration maintains a *compact* value-policy pair  $(r_t, \pi_t)$ . Given  $r_t$ ,  $\pi_{t+1}$  is the greedy policy with respect to

---

18. Bertsekas and Ioffe (1996) consider a weaker assumption for Exact  $\lambda$  Policy Iteration and its analysis, namely that there exists at least *one* proper policy. However, and as will be clear by its description, this assumption is not sufficient for their Approximate algorithm, because this builds sample trajectories that need to reach a termination state. If the terminal state were not reachable in finite time, this algorithm may not terminate in finite time.

$v^{r_t}$ , and can easily be computed exactly in any given state as the argmax in Equation 19. This policy  $\pi_{t+1}$  is used to simulate a batch of  $M$  trajectories: for each trajectory  $m$ ,  $(s_{m,0}, s_{m,1}, \dots, s_{m,N_m-1}, s_{m,N_m})$  denotes the sequence of states of the  $m$ th trajectory, with  $s_{m,N_m}$  being the termination state. Then for approximating Equation 13, a reasonable choice for  $r_{t+1}$  is one that satisfies:

$$\begin{aligned}
v^{\mathbf{I}^{t+1}}(s_{m,N_m}) &\simeq 0 \\
v^{\mathbf{I}^{t+1}}(s_{m,N_m-1}) &\simeq v^{r_t}(s_{m,N_m-1}) + \delta_t(s_{m,N_m-1}, s_{m,N_m}) \\
v^{\mathbf{I}^{t+1}}(s_{m,N_m-2}) &\simeq v^{r_t}(s_{m,N_m-2}) + \delta_t(s_{m,N_m-2}, s_{m,N_m-1}) + \gamma\lambda\delta_t(s_{m,N_m-1}, s_{m,N_m}) \\
&\vdots \\
v^{\mathbf{I}^{t+1}}(s_{m,k}) &\simeq v^{r_t}(s_{m,k}) + \sum_{s=k}^{N_m-1} (\gamma\lambda)^{s-k} \delta_t(s_{m,s}, s_{m,s+1}) \\
&\vdots \\
v^{\mathbf{I}^{t+1}}(s_{m,0}) &\simeq v^{r_t}(s_{m,0}) + \sum_{s=0}^{N_m-1} (\gamma\lambda)^s \delta_t(s_{m,s}, s_{m,s+1})
\end{aligned} \tag{22}$$

for all trajectories  $m$ , where

$$\delta_t(s_{m,N_m-1}, s_{m,N_m}) = r(s_{m,N_m-1}, \pi_{t+1}(s_{m,N_m-1}), s_{m,N_m}) - v^{r_t}(s_{m,N_m-1}) \tag{23}$$

and for all  $s < N_m - 1$

$$\delta_t(s_{m,s}, s_{m,s+1}) = r(s_{m,s}, \pi_{t+1}(s_{m,s}), s_{m,s+1}) + \gamma v^{r_t}(s_{m,s+1}) - v^{r_t}(s_{m,s})$$

are the temporal differences. Note well that Equations 22 and 23 correspond to the terminal states for which there is no subsequent rewards. A standard and efficient solution to this problem consists in minimizing the least squares error, that is to choose  $r_{t+1}$  as follows:

$$r_{t+1} = \arg \min_{\mathbf{r}} \sum_{m=1}^M \sum_{k=0}^{N_m} \left( v^{\mathbf{r}}(s_{m,k}) - v^{r_t}(s_{m,k}) - \sum_{s=k}^{N_m-1} (\gamma\lambda)^{s-k} \delta_t(s_{m,s}, s_{m,s+1}) \right)^2.$$

This approximate version of  $\lambda$  Policy Iteration generalizes well known algorithms. When  $\lambda = 0$ , the generic term becomes a sample of  $[\mathcal{T}^{\pi_{k+1}}v](s_{m,k})$ :

$$v^{\mathbf{I}^{t+1}}(s_{m,k}) \simeq v^{r_t}(s_{m,k}) + \delta_t(s_{m,k}, s_{m,k+1}) = r(s_{m,k}, \pi_{t+1}(s_{m,k}), s_{m,k+1}) + \gamma v^{r_t}(s_{m,k+1}) \tag{24}$$

When  $\lambda = 1$ , the generic term becomes the sampled discounted return from  $s_{m,k}$  until the end of the trajectory:

$$v^{\mathbf{I}^{t+1}}(s_{m,k}) \simeq v^{r_t}(s_{m,k}) + \sum_{s=k}^{N_m-1} \gamma^{s-k} \delta_t(s_{m,s}, s_{m,s+1}) = \sum_{s=k}^{N_m-1} \gamma^{s-k} r(s_{m,k}, \pi_{t+1}(s_{m,k}), s_{m,k+1}) \tag{25}$$

In other words, for these limit values of  $\lambda$ , the algorithms correspond to approximate versions of Value Iteration and Policy Iteration as described by Bertsekas and Tsitsiklis (1996).

Also, as explained in Bertsekas and Ioffe (1996) and already mentioned in the introduction, the TD( $\lambda$ ) algorithm with linear features described in chapter 8.2 by Sutton and Barto (1998) matches the algorithm we have just described when the above fitting problem is *approximated* using gradient iterations after each sample.

We follow the same protocol as originally proposed by Bertsekas and Ioffe (1996). Let  $w = 10$  be the width of the board. We consider approximating the value function as a linear combination of  $2w + 1 = 21$  feature functions:

$$V^r(s) = r(0) + \sum_{k=1}^w r(k)h_k + \sum_{k=1}^{w-1} r(k+w)\Delta h_k + r(2w)H + r(2w+1)L$$

where:

- for all  $k \in \{1, 2, \dots, w\}$ ,  $h_k$  is the *height* of the  $k$ th column of the wall;
- For all  $k \in \{1, 2, \dots, w-1\}$ ,  $\Delta h_k$  is the *height difference*  $|h_k - h_{k+1}|$  between columns  $k$  and  $k+1$ ;
- $H$  is the *maximum wall height*  $\max_k h_k$ ;
- $L$  is the number of *holes* (the number of empty cells recovered by at least one full cell).

We started our experiments with the initial following vector:  $r(2w) = -10$ ,  $r(2w+1) = -1$  and  $r(k) = 0$  for all  $k < 2w$ , so that the initial greedy policy scores in the low tens (Bertsekas and Ioffe (1996)). We used  $M = 100$  training games for each policy update. As  $\lambda PI$  is a stochastic algorithm, we ran each experiment 10 times. Figure 4 displays the

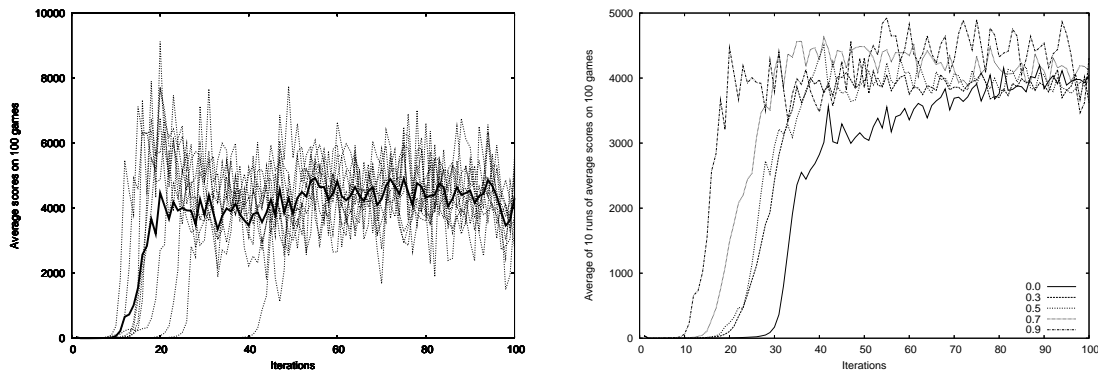


Figure 4: Average Score versus the number of iterations **Left**: 10 runs of  $\lambda PI$  with  $\lambda = 0.9$ . Each point of each run is the average score computed with  $M = 100$  games. The dark curve is a pointwise average of the 10 runs. **Right**: Pointwise average of 10 runs of  $\lambda PI$  for different values of  $\lambda$ ; the curve which appears to be the best ( $\lambda = 0.9$ ) is the same as the bold curve of the left graph.

learning curves. The left graph shows the 10 runs of  $\lambda PI$  (each point is the average score computed with the  $M = 100$  games) and the corresponding pointwise average for a single



value of  $\lambda$ , while the right graph shows such pointwise average curves for different values of  $\lambda$ : 0.0, 0.3, 0.5, 0.7 and 0.9. We chose to display on the left graph the runs corresponding to the value of  $\lambda = 0.9$  that seemed to be the best on the right graph.

We can make the following observations:

- Though we initialized with not so bad a policy (the first value is around 30), the performance first drops to 0 and it *really* starts improving after a few iterations (typically around ten). This is due to the fact that the initial value function is really bad: with the given parameters, the initial value is everywhere negative although it is clear that the optimal value function (the average best score) is everywhere positive. Further experiments showed that the overall behaviour of the algorithm was not affected by the weight initialization.
- The rise of performance globally happens sooner for larger values of  $\lambda$ , that is for values that makes the algorithm closer to Policy Iteration. This is not surprising as it complies with the fact that  $\lambda$  modulates the speed at which the value estimate tracks the real value of the current policy. However, the performance did not rise for  $\lambda = 1$  (when it is equivalent to Approximate Policy Iteration), and this is probably due to the fact that the variance of the value update is too high.
- Quantitatively, the scores reach an overall level of 4000 lines per games for a big range of values of  $\lambda$ .

The empirical results we have just described qualitatively and quantitatively differ from the ones that were originally published in Bertsekas and Ioffe (1996), even though it is the exact same experimental setup. About their results, the authors wrote: “*An interesting and somewhat paradoxical observation is that a high performance is achieved after relatively few policy iterations, but the performance gradually drops significantly. We have no explanation for this intriguing phenomenon, which occurred with all of the successful methods that we tried*”. As we explain now, we believe that the “intriguing” character of the results of Bertsekas and Ioffe (1996) might be related to a subtle implementation difference. Indeed, we can reproduce learning curves that are similar to those of Bertsekas and Ioffe (1996) with a little modification in our implementation of  $\lambda PI$ , that removes the special treatments for the terminal states done through Equations 22 and 23. More precisely, if we replace them by the following Equations:

$$v^{\mathbf{I}^{t+1}}(s_{m,N_m}) \simeq v^{rt}(s_{m,N_m}) \quad (26)$$

$$\delta_t(s_{m,N_m-1}, s_{m,N_m}) = r(s_{m,N_m-1}, \pi_{t+1}(s_{m,N_m-1})) + \gamma V^{rt}(s_{m,N_m}) - V^{rt}(s_{m,N_m-1}) \quad (27)$$

that is if we replace the terminal value 0 by the value  $V^{rt}(s_{m,N_m})$  which is computed through the features of the terminal wall configuration  $s_{m,N_m}$ , then we get the performance shown in figure 5. This figure shows the performance with respect to the iterations. We observe that the performance evolution qualitatively matches the performance curves published in Bertsekas and Ioffe (1996) and illustrates the above quotation describing the “intriguing phenomenon”<sup>19</sup>.

---

19. The watchful reader may have noticed that the performance that we obtain is about twice that of Bertsekas and Ioffe (1996). A close inspection of the Tetris domain description in (Bertsekas and Ioffe,

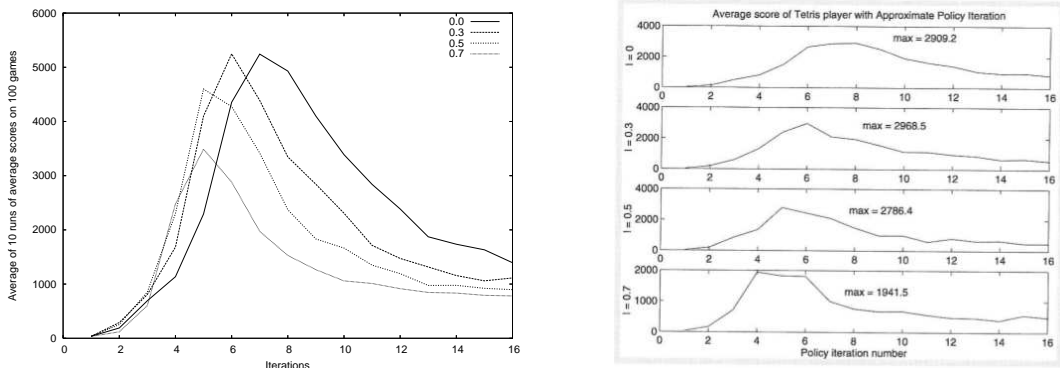


Figure 5: **Left:** Average score versus the number of iterations of  $\lambda PI$ , *modified* so that it resembles the results of Bertsekas and Ioffe (1996) (see text for details). **Right:** A scan of the learning curves obtained by Bertsekas and Ioffe (1996).

In such a modified form, the approximate  $\lambda PI$  algorithm makes much less sense: in particular, it is not true anymore that it reduces to approximate Value Iteration and approximate Policy Iteration when  $\lambda = 0$  and  $\lambda = 1$  respectively: Equations 26 and 27 induces a bias so that we cannot recover the identities of Equations 24 and 25. A closer examination of these experiments showed that the weights  $(r_k)$  were diverging. This is not a surprise, since the use of Equations 26 and 27 violates the condition expressed in Remark 31-1 that there should be no error in the terminal state.

## 8. Conclusion and Future Work

We have considered the  $\lambda$  Policy Iteration algorithm introduced by Bertsekas and Ioffe (1996) that generalizes the standard algorithms Value Iteration and Policy Iteration. We have reviewed some results by Puterman (1994), Bertsekas and Tsitsiklis (1996) and Munos (2003, 2007), concerning the rate of convergence and some performance bounds of these standard algorithms in the exact and approximate cases. We have extended these results to  $\lambda$  Policy Iteration and derived convergence rates and performance bound for this algorithm, some of which are to our knowledge even new in the cases where  $\lambda$  Policy Iteration reduces to Value Iteration and Policy Iteration, notably the bound 16 (page 20). Not only does our analysis generalize previous results, but it improves them in two ways:

- As suggested by the results of Puterman (1994) in the exact case, the use of the span seminorm has enabled us to derive tighter bounds.

---

1996) shows that the authors consider the game of Tetris on a  $10 \times 19$  board instead of our  $10 \times 20$  setting, and as argued in a recent review on Tetris (Thiery and Scherrer, 2009), this small difference is sufficient for explaining such a big performance difference.

- Our analysis of Approximate  $\lambda$  Policy Iteration relates the asymptotic performance of the algorithm to the *asymptotic* errors/residuals instead of a uniform bound of the errors/residuals and this might be of practical interest<sup>20</sup>.

More generally, we believe that an important contribution of this paper is of conceptual nature: we provide a unified vision of some of the main Approximate Dynamic Programming algorithms and their analyses; in particular, we hope that the new proof technique that is detailed in the appendices (especially the different objects that are defined in our proof overview (Section 4.2)) will be useful for further works.

As we mentioned earlier, Munos (2007) introduced some concentration coefficients that are finer than the one we used throughout the paper. In the same spirit, Farahmand *et al.* (2010) recently revisited the error propagation of Munos (2007, 2003) and improved (among other things) the constant in the bound related to these concentration coefficients. A natural track would be to adapt these refinements to the analysis  $\lambda$  Policy Iteration. This does not look completely trivial since the componentwise analysis we derived for  $\lambda$  Policy Iteration is significantly more intricate than the ones we find in the specific limit cases  $\lambda = 0$  (Value Iteration) and  $\lambda = 1$  (Policy Iteration).

Another potential direction would be to study the implications of the choice of the parameter  $\lambda$ , as for instance is done by Singh and Dayan (1998) for the value estimation problem. On this matter, the original analysis by Bertsekas and Ioffe (1996) shows how one can concretely implement Exact  $\lambda$  Policy Iteration. Each iteration requires the computation of the fixed point of the  $\beta$ -contracting operator  $M_k$  (see Equation 9 page 11). We plan to study the tradeoff between the ease for computing this fixed point (the smaller  $\beta$  the faster) and the time for  $\lambda$  Policy Iteration to converge to the optimal policy (the bigger  $\beta$  the faster). In parallel, the reader might have noticed that most of the bounds we have provided do not depend on  $\lambda$ . An interesting question is whether the finer concentration coefficients of Munos (2007) and Farahmand *et al.* (2010) we have just discussed may help keeping track of the influence of  $\lambda$  on the performance of the exact or approximate algorithm. In general, it would be interesting if we could tie a choice of  $\lambda$  to some intrinsic characteristics of the MDP, like for instance its smoothness.

Last but not least, we should insist on the fact that the implementation that we have described in Section 7.2, and which is borrowed from Bertsekas and Ioffe (1996), is just one possible instance of  $\lambda$  Policy Iteration. In the case of linear approximation architectures, Thiery and Scherrer (2010) have proposed an implementation of  $\lambda$  Policy Iteration that is based on LSPI (Lagoudakis and Parr, 2003), in which the fixed point of  $M_k$  is approximated using LSTD(0) (Bradtke and Barto, 1996). Recently, Bertsekas (2011) proposed to compute this very fixed point with a variation of LSPE( $\lambda'$ ) (Bertsekas and Ioffe, 1996; Nedić and Bertsekas, 2003) for some  $\lambda'$  potentially different from  $\lambda$ . Because of their very close structure, any existing implementation of Approximate Policy Iteration may probably be turned into some implementation of  $\lambda$  Policy Iteration. Proposing such implementations and assessing their relative merits constitutes interesting future research. This may in particular be done through some finite sample analysis, as had recently been done for Approximate

---

20. Recently and independently, Farahmand *et al.* (2010) derived bounds that have a similar flavour: they highlight the fact that the errors that have the more weight on the performance bounds are the latest.

Value Iteration and Policy Iteration implementations (Antos *et al.*, 2007, 2008; Munos and Szepesvári, 2008; Lazaric *et al.*, 2010).

## References

- Antos, A., Szepesvári, C., and Munos, R. (2007). Value-iteration based fitted policy iteration: Learning with a single trajectory. In *ADPRL 2007*, pages 330–337. IEEE.
- Antos, A., Szepesvári, C., and Munos, R. (2008). Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning Journal*, **71**, 89–129.
- Bertsekas, D. (2011). Lambda Policy Iteration: A Review and A New Implementation. Technical Report LIDS-2874, MIT.
- Bertsekas, D. and Ioffe, S. (1996). Temporal differences-based policy iteration and applications in neuro-dynamic programming. Technical Report LIDS-P-2349, MIT.
- Bertsekas, D. and Tsitsiklis, J. (1996). *Neurodynamic Programming*. Athena Scientific.
- Bradtke, S. J. and Barto, A. G. (1996). Linear Least-Squares algorithms for temporal difference learning. *Machine Learning*, **22**(1-3), 33–57.
- Burgiel, H. (1997). How to Lose at Tetris. *Mathematical Gazette*, **81**, 194–200.
- Fahey, C. P. (2003). Tetris AI, Computer plays Tetris. [http://colinfahey.com/tetris/tetris\\_en.html](http://colinfahey.com/tetris/tetris_en.html).
- Farahmand, A., Munos, R., and Szepesvári, C. (2010). Error Propagation for Approximate Policy and Value Iteration. In *NIPS*.
- Gordon, G. (1995). Stable function approximation in dynamic programming. In A. Prieditis and S. Russell, editors, *ICML*, pages 261–268, San Francisco, CA. Morgan Kaufmann.
- Lagoudakis, M. and Parr, R. (2003). Least-Squares Policy Iteration. *Journal of Machine Learning Research*, **4**, 1107–1149.
- Lazaric, A., Ghavamzadeh, M., and Munos, R. (2010). Analysis of a Classification-based Policy Iteration Algorithm. In *ICML*, pages 607–614.
- Munos, R. (2003). Error Bounds for Approximate Policy Iteration. In *ICML*, pages 560–567.
- Munos, R. (2007). Performance bounds in  $L_p$  norm for approximate value iteration. *SIAM J. Control and Optimization*.
- Munos, R. and Szepesvári, C. (2008). Finite-Time Bounds for Fitted Value Iteration. *Journal of Machine Learning Research*, **9**, 815–857.
- Nedić, A. and Bertsekas, D. P. (2003). Least Squares Policy Evaluation Algorithms with Linear Function Approximation. *DEDS*, **13**, 79–110.

- Puterman, M. (1994). *Markov Decision Processes*. Wiley, New York.
- Puterman, M. and Shin, M. (1978). Modified policy iteration algorithms for discounted Markov decision problems. *Management Science*, **24**(11).
- Singh, S. and Dayan, P. (1998). Analytical Mean Squared Error Curves for Temporal Difference Learning. *Machine Learning Journal*, **32**(1), 5–40.
- Sutton, R. and Barto, A. (1998). *Reinforcement Learning, An introduction*. Bradford Book. The MIT Press.
- Thiery, C. and Scherrer, B. (2009). Improvements on Learning Tetris with Cross Entropy. *International Computer Games Association Journal*, **32**.
- Thiery, C. and Scherrer, B. (2010). Least-Squares  $\lambda$  Policy Iteration: Bias-Variance Trade-off in Control Problems. In *ICML*, Haifa, Israël.
- Van Roy, B. (2006). Performance Loss Bounds for Approximate Value Iteration with State Aggregation. *Mathematics of Operation Research*, **31**(2), 234–244.
- Williams, R. and Baird, L. (1993). Tight Performance Bounds on Greedy Policies Based on Imperfect Value Functions. Technical Report NU-CCS-93-14, Northeastern University.

## Appendices

The following Appendices contains all the proofs concerning the analysis of  $\lambda$  Policy Iteration. We write  $P_k = P^{\pi_k}$  the stochastic matrix corresponding to the policy  $\pi_k$  which is greedy with respect to  $v_{k-1}$ ,  $P_*$  the stochastic matrix corresponding to the optimal policy  $\pi_*$ . Similarly we write  $\mathcal{T}_k$  and  $\mathcal{T}_*$  the associated Bellman operators.

The proof techniques we have developped are inspired by those of Munos in the articles Munos (2003, 2007). Most of the inequalities will appear from the definition of the greedy operator:

$$\pi = \text{greedy}(v) \Leftrightarrow \forall \pi', \mathcal{T}^{\pi'} v \leq \mathcal{T}^{\pi} v.$$

We will often use the property that an average of stochastic matrices is also a stochastic matrix. A recurrent instance of this property is: if  $P$  is some stochastic matrix, then the geometric average

$$(1 - \alpha) \sum_{i=0}^{\infty} (\alpha P)^i = (1 - \alpha)(I - \alpha P)^{-1}$$

with  $0 \leq \alpha < 1$  is also a stochastic matrix. Finally, we will use the property that if some vectors  $x$  and  $y$  are such that  $x \leq y$ , then  $Px \leq Py$  for any stochastic matrix  $P$ .

## Appendix A. General componentwise analysis of $\lambda$ Policy Iteration

In this section, we prove the series of Lemmas that are useful for understanding the error propagation of  $\lambda$  Policy Iteration.

### A.1 Proof of Lemma 13: a relation between the shift and the Bellman residual

We have:

$$\begin{aligned}
 (I - \gamma P_k)s_k &= (I - \gamma P_k)(w_k - v^{\pi_k}) \\
 &= (I - \gamma P_k)w_k - r_k \\
 &= (I - \lambda\gamma P_k + \lambda\gamma P_k - \gamma P_k)w_k - r_k \\
 &= (I - \lambda\gamma P_k)w_k + (\lambda\gamma P_k - \gamma P_k)w_k - r_k \\
 &= r_k + (1 - \lambda)\gamma P_k v_{k-1} + (\lambda - 1)\gamma P_k w_k - r_k \\
 &= (1 - \lambda)\gamma P_k(v_{k-1} - w_k) \\
 &= (1 - \lambda)\gamma P_k(I - \lambda\gamma P_k)^{-1}(v_{k-1} - \mathcal{T}_k v_{k-1}) \\
 &= (1 - \lambda)\gamma P_k(I - \lambda\gamma P_k)^{-1}(-b_{k-1})
 \end{aligned}$$

Therefore

$$s_k = \beta(I - \gamma P_k)^{-1}A_k(-b_{k-1})$$

with

$$A_k := (1 - \lambda\gamma)P_k(I - \lambda\gamma P_k)^{-1}.$$

Suppose we have a lower bound of the Bellman residual:  $b_k \geq \underline{b}_k$  (we will derive one soon). Since  $(I - \gamma P_k)^{-1}A_k$  only has non-negative elements then

$$s_k \leq \beta(I - \gamma P_k)^{-1}A_k(-\underline{b}_k) := \overline{s}_k.$$

### A.2 Proof of Lemma 14: a lower bound of the Bellman residual

From the definition of the algorithm, and using the fact that  $\mathcal{T}_k v^{\pi_k} = v^{\pi_k}$  we see that:

$$\begin{aligned}
 b_k &= \mathcal{T}_{k+1}v_k - v_k \\
 &= \mathcal{T}_{k+1}v_k - \mathcal{T}_k v_k + \mathcal{T}_k v_k - v_k \\
 &\geq \mathcal{T}_k v_k - v_k \\
 &= \mathcal{T}_k v_k - \mathcal{T}_k v^{\pi_k} + v^{\pi_k} - v_k \\
 &= \gamma P_k(v_k - v^{\pi_k}) + v^{\pi_k} - v_k \\
 &= (\gamma P_k - I)(s_k + \epsilon_k). \\
 &= \beta A_k b_{k-1} + (\gamma P_k - I)\epsilon_k.
 \end{aligned} \tag{28}$$

where we eventually used the relation between  $s_k$  and  $b_k$  (Lemma 13). In other words:

$$b_{k+1} \geq \beta A_{k+1} b_k + x_{k+1}$$

with

$$x_k := (\gamma P_k - I)\epsilon_k.$$

Since  $A_k$  is a stochastic matrix and  $\beta \geq 0$ , we get by induction:

$$b_k \geq \sum_{j=k_0+1}^k \beta^{k-j} (A_k A_{k-1} \dots A_{j+1}) x_j + \beta^{k-k_0} (A_k A_{k-1} \dots A_{k_0+1}) b_{k_0} := \underline{b}_k.$$

### A.3 Proof of Lemma 15: an upper bound of the distance

Given that  $\mathcal{T}_* v_* = v_*$ , we have

$$\begin{aligned} v_* &= v_* + (I - \lambda\gamma P_{k+1})^{-1} (\mathcal{T}_* v_* - v_*) \\ &= (I - \lambda\gamma P_{k+1})^{-1} (\mathcal{T}_* v_* - \lambda\gamma P_{k+1} v_*) \end{aligned}$$

Therefore the distance satisfies:

$$\begin{aligned} d_{k+1} &= v_* - w_{k+1} \\ &= (I - \lambda\gamma P_{k+1})^{-1} [(\mathcal{T}_* v_* - \lambda\gamma P_{k+1} v_*) - (\mathcal{T}_{k+1} v_k - \lambda\gamma P_{k+1} v_k)] \\ &= (I - \lambda\gamma P_{k+1})^{-1} [\mathcal{T}_* v_* - \mathcal{T}_{k+1} v_k + \lambda\gamma P_{k+1} (v_k - v_*)] \\ &= \lambda\gamma P_{k+1} d_{k+1} + \mathcal{T}_* v_* - \mathcal{T}_{k+1} v_k + \lambda\gamma P_{k+1} (v_k - v_*) \\ &= \lambda\gamma P_{k+1} d_{k+1} + \mathcal{T}_* v_* - \mathcal{T}_{k+1} v_k + \lambda\gamma P_{k+1} (w_k + \epsilon_k - v_*) \\ &= \lambda\gamma P_{k+1} d_{k+1} + \mathcal{T}_* v_* - \mathcal{T}_{k+1} v_k + \lambda\gamma P_{k+1} (\epsilon_k - d_k) \\ &= \mathcal{T}_* v_* - \mathcal{T}_{k+1} v_k + \lambda\gamma P_{k+1} (\epsilon_k + d_{k+1} - d_k) \end{aligned}$$

Since  $\pi^{k+1}$  is greedy with respect to  $v_k$ , we have  $\mathcal{T}_{k+1} v_k \geq \mathcal{T}_* v_k$  and therefore:

$$\begin{aligned} \mathcal{T}_* v_* - \mathcal{T}_{k+1} v_k &= \mathcal{T}_* v_* - \mathcal{T}_* v_k + \mathcal{T}_* v_k - \mathcal{T}_{k+1} v_k \\ &\leq \mathcal{T}_* v_* - \mathcal{T}_* v_k \\ &= \gamma P_* (v_* - v_k) \\ &= \gamma P_* (v_* - (w_k + \epsilon_k)) \\ &= \gamma P_* d_k - \gamma P_* \epsilon_k \end{aligned}$$

As a consequence, the distance satisfies:

$$d_{k+1} \leq \gamma P_* d_k + \lambda\gamma P_{k+1} (\epsilon_k + d_{k+1} - d_k) - \gamma P_* \epsilon_k$$

Noticing that:

$$\begin{aligned} \epsilon_k + d_{k+1} - d_k &= \epsilon_k + w_k - w_{k+1} \\ &= v_k - w_{k+1} \\ &= -(I - \lambda\gamma P_{k+1})^{-1} (\mathcal{T}_{k+1} v_k - v_k) \\ &= (I - \lambda\gamma P_{k+1})^{-1} (-b_k) \\ &\leq (I - \lambda\gamma P_{k+1})^{-1} (-\underline{b}_k) \end{aligned}$$

we get:

$$d_{k+1} \leq \gamma P_* d_k + y_k$$

where

$$y_k := \frac{\lambda\gamma}{1-\lambda\gamma} A_{k+1}(-\underline{b}_k) - \gamma P_* \epsilon_k.$$

Since  $P_*$  is a stochastic matrix and  $\gamma \geq 0$ , we have by induction:

$$d_k \leq \sum_{j=k_0}^{k-1} \gamma^{k-1-j} (P_*)^{k-1-j} y_j + \gamma^{k-k_0} (P_*)^{k-k_0} d_{k_0} = \overline{d}_k.$$

## Appendix B. Componentwise rate of convergence for Exact $\lambda$ Policy Iteration

We here derive the convergence rate bounds for Exact  $\lambda$  Policy Iteration (as expressed in Lemma 17 page 19). We rely on the loss bound analysis of Appendix A with  $\epsilon_k = 0$ . In this specific case, we know that the loss  $l_k \leq \overline{d}_k + \overline{s}_k$  where

$$\begin{aligned} -\underline{b}_k &= \beta^{k-k_0} A_k A_{k-1} \dots A_{k_0+1} (-b_{k_0}) \\ \overline{d}_k &= \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{j=k_0}^{k-1} \gamma^{k-1-j} (P_*)^{k-1-j} A_{j+1} (-\underline{b}_j) + \gamma^{k-k_0} (P_*)^{k-k_0} d_{k_0} \\ \overline{s}_k &= \beta (I - \gamma P_k)^{-1} A_k (-\underline{b}_{k-1}). \end{aligned}$$

Introducing the following stochastic matrices:

$$X_{i,j,k} := (P_*)^{k-1-i} A_{i+1} A_i \dots A_{j+1}$$

$$Y_{j,k} := (1-\gamma)(I - \gamma P_k)^{-1} A_k A_{k-1} \dots A_{j+1},$$

we have

$$\overline{d}_k = \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{j=k_0}^{k-1} \gamma^{k-1-j} \beta^{j-k_0} X_{j,k_0,k} (-b_{k_0}) + \gamma^{k-k_0} (P_*)^{k-k_0} d_{k_0}$$

and

$$\overline{s}_k = \frac{\beta^{k-k_0}}{1-\gamma} Y_{k_0,k} (-b_{k_0}).$$

Therefore

$$l_k \leq \left( \frac{\gamma^{k-k_0}}{1-\gamma} \right) E'_{kk_0} (-b_{k_0}) + \gamma^{k-k_0} (P_*)^{k-k_0} d_{k_0} \quad (29)$$

with

$$E'_{kk_0} := \left( \frac{1-\gamma}{\gamma^{k-k_0}} \right) \left( \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{j=k_0}^{k-1} \gamma^{k-1-j} \beta^{j-k_0} X_{j,k_0,k} + \frac{\beta^{k-k_0}}{1-\gamma} Y_{k_0,k} \right).$$

**Lemma 32**  $E'_{kk_0}$  is a stochastic matrix



**Proof** It is clear from the definition of  $X_{i,j,k}$  and  $Y_{j,k}$  that normalizing  $E'_{kk_0}$  will give stochastic matrices. So we just need to check that their max-norm is 1.

$$\begin{aligned}
\|E'_{kk_0}\| &= \frac{1-\gamma}{\gamma^{k-k_0}} \left( \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{j=k_0}^{k-1} \gamma^{k-1-j} \beta^{j-k_0} + \frac{\beta^{k-k_0}}{1-\gamma} \right) \\
&= \frac{1-\gamma}{\gamma^{k-k_0}} \left( \frac{\lambda\gamma}{1-\lambda\gamma} \frac{\gamma^{k-k_0} - \beta^{k-k_0}}{\gamma - \beta} + \frac{\beta^{k-k_0}}{1-\gamma} \right) \\
&= \frac{1-\gamma}{\gamma^{k-k_0}} \left( \frac{\gamma^{k-k_0} - \beta^{k-k_0}}{1-\gamma} + \frac{\beta^{k-k_0}}{1-\gamma} \right) \\
&= 1
\end{aligned}$$

where we used the facts that  $\frac{\lambda\gamma}{\gamma-\beta} = \frac{1}{1-\beta}$  and  $(1-\beta)(1-\lambda\gamma) = 1-\gamma$ . ■

### B.1 A bound with respect to the Bellman residual

We first need the following lemma:

**Lemma 33** *The bias and the distance are related as follows:*

$$b_k \geq (I - \gamma P_*) d_k.$$

**Proof** Since  $\pi_{k+1}$  is greedy with respect to  $v_k$ ,  $\mathcal{T}_{k+1}v_k \geq \mathcal{T}_*v_k$  and

$$\begin{aligned}
b_k &= \mathcal{T}_{k+1}v_k - v_k \\
&= \mathcal{T}_{k+1}v_k - \mathcal{T}_*v_k + \mathcal{T}_*v_k - \mathcal{T}_*v_* + v_* - v_k \\
&\geq \gamma P_*(v_k - v_*) + v_* - v_k \\
&= (I - \gamma P_*) d_k.
\end{aligned}$$
■

We thus have:

$$d_{k_0} \leq (I - \gamma P_*)^{-1} b_{k_0}.$$

Then Equation 29 becomes

$$\begin{aligned}
l_k &\leq \left[ \gamma^{k-k_0} (P_*)^{k-k_0} (I - \gamma P_*)^{-1} - \left( \frac{\gamma^{k-k_0}}{1-\gamma} \right) E'_{kk_0} \right] b_{k_0} \\
&= \frac{\gamma^{k-k_0}}{1-\gamma} [E_{kk_0} - E'_{kk_0}] b_{k_0}
\end{aligned}$$

where:

$$E_{kk_0} := (1-\gamma)(P_*)^{k-k_0}(I - \gamma P_*)^{-1}$$

is a stochastic matrix.

## B.2 A bound with respect to the distance

From Lemma 33, we know that

$$-b_{k_0} \leq (I - \gamma P_*)(-d_{k_0})$$

Then Equation 29 becomes

$$\begin{aligned} l_k &\leq \left[ \gamma^{k-k_0} (P_*)^{k-k_0} - \left( \frac{\gamma^{k-k_0}}{1-\gamma} \right) E'_{kk_0} (I - \gamma P_*) \right] d_{k_0} \\ &= \frac{\gamma^{k-k_0}}{1-\gamma} [F_{kk_0} - E'_{kk_0}] d_{k_0} \end{aligned}$$

where

$$F_{kk_0} := (1-\gamma)P_*^{k-k_0} + \gamma E'_{kk_0} P_*$$

is a stochastic matrix.

## B.3 A bound with respect to the distance and the loss of the greedy policy

Let  $K$  be a constant such that  $\hat{v}_{k_0} := v_{k_0} - Ke$ . The following statements are equivalent:

$$\begin{aligned} \hat{b}_{k_0} &\geq 0 \\ \mathcal{T}_{k_0+1} \hat{v}_{k_0} &\geq \hat{v}_{k_0} \\ r_{k_0+1} + \gamma P_{k_0+1} (v_{k_0} - Ke) &\geq v_{k_0} - Ke \\ (I - \gamma P_{k_0+1}) Ke &\geq -r_{k_0+1} + (I - \gamma P_{k_0+1}) v_{k_0} \\ Ke &\geq (I - \gamma P_{k_0+1})^{-1} (-r_{k_0+1}) + v_{k_0} \\ Ke &\geq v_{k_0} - v^{\pi_{k_0+1}} \end{aligned}$$

The minimal  $K$  for which  $\hat{b}_{k_0} \geq 0$  is thus  $K := \max_s [v_{k_0}(s) - v^{\pi_{k_0+1}}(s)]$ . As  $\hat{v}_{k_0}$  and  $v_{k_0}$  only differ by a constant vector, they will generate the same sequence of policies  $\pi_{k_0+1}, \pi_{k_0+2}, \dots$ . Then, as  $\hat{b}_{k_0} \geq 0$ , Equation 29 tells us that

$$\begin{aligned} v_* - v^{\pi_k} &\leq \gamma^{k-k_0} (P_*)^{k-k_0} (v_* - \hat{v}_{k_0}) \\ &= \gamma^{k-k_0} (P_*)^{k-k_0} (v_* - v_{k_0} + Ke). \end{aligned}$$

Now notice that

$$\begin{aligned} K &= \max_s [v_{k_0}(s) - v_*(s) + v_*(s) - v^{\pi_{k_0+1}}(s)] \\ &\leq \max_s [v_{k_0}(s) - v_*(s)] + \max_s [v_*(s) - v^{\pi_{k_0+1}}(s)] \\ &= -\min_s [v_*(s) - v_{k_0}(s)] + \|v_*(s) - v^{\pi_{k_0+1}}\|_\infty \end{aligned}$$

Then, using the fact that  $(P_*)^{k-k_0} e = e$ , we get:

$$v_* - v^{\pi_k} \leq \gamma^{k-k_0} \left[ (P_*)^{k-k_0} \left( (v_* - v_{k_0}) - \min_s [v_*(s) - v_{k_0}(s)] e \right) + \|v_*(s) - v^{\pi_{k_0+1}}\|_\infty e \right].$$

## Appendix C. Asymptotic componentwise loss bounds with respect to the approximation error

We here use the loss bound analysis of Appendix A to derive an asymptotic analysis of approximate  $\lambda$  Policy Iteration with respect to the approximation error. The results stated here constitute a proof of the first inequality of Lemma 21 page 21, and of Corollary 27 page 25.

### C.1 General analysis

Since

$$l_k = d_k + s_k \leq \overline{d}_k + \overline{s}_k, \quad (30)$$

an upper bound of the loss can be derived from the upper bound of the distance and the shift.

Let us first concentrate on the bound  $\overline{d}_k$  of the distance. So far we have proved that:

$$\begin{aligned} \overline{d}_k &= \sum_{i=k_0}^{k-1} \gamma^{k-1-i} (P_*)^{k-1-i} y_i + \mathcal{O}(\gamma^{k-k_0}) \\ y_i &= \frac{\lambda\gamma}{1-\lambda\gamma} A_{i+1}(-\underline{b}_i) - \gamma P_* \epsilon_i \\ -\underline{b}_i &= \sum_{j=k_0}^i \beta^{i-j} (A_i A_{i-1} \dots A_{j+1}) (-x_j) + \mathcal{O}(\beta^{i-k_0}) \\ -x_j &= (I - \gamma P_j) \epsilon_j. \end{aligned}$$

Writing

$$X_{i,j,k} := (P_*)^{k-1-i} A_{i+1} A_i \dots A_{j+1}$$

and putting all things together, we see that:

$$\begin{aligned} \overline{d}_k &= \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=k_0}^{k-1} \gamma^{k-1-i} \left( \sum_{j=k_0}^i \beta^{i-j} X_{i,j,k} (I - \gamma P_j) \epsilon_j + \mathcal{O}(\beta^{i-k_0}) \right) \\ &\quad - \sum_{i=k_0}^{k-1} \gamma^{k-i} (P_*)^{k-i} \epsilon_i + \mathcal{O}(\gamma^{k-k_0}) \\ &= \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=k_0}^{k-1} \sum_{j=k_0}^i \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} (I - \gamma P_j) \epsilon_j - \sum_{i=k_0}^{k-1} \gamma^{k-i} (P_*)^{k-i} \epsilon_i + \mathcal{O}(\gamma^{k-k_0}) \\ &= \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{j=k_0}^{k-1} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} (I - \gamma P_j) \epsilon_j - \sum_{j=k_0}^{k-1} \gamma^{k-j} (P_*)^{k-j} \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \\ &= \sum_{j=k_0}^{k-1} \left[ \left( \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} (I - \gamma P_j) \right) - \gamma^{k-j} (P_*)^{k-j} \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \quad (31) \end{aligned}$$

where between the first two lines, we used the fact that:

$$\frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=k_0}^{k-1} \gamma^{k-1-i} \beta^{i-k_0} = \frac{\lambda\gamma}{1-\lambda\gamma} \frac{\gamma^{k-k_0} - \beta^{k-k_0}}{\gamma - \beta} = \frac{\gamma^{k-k_0} - \beta^{k-k_0}}{1-\gamma} = \mathcal{O}(\gamma^{k-k_0}) \quad (32)$$

using the identities  $\lambda\gamma = \frac{\gamma-\beta}{1-\beta}$  and  $1-\gamma\lambda = \frac{1-\gamma}{1-\beta}$ .

Let us now consider the bound  $\overline{s}_k$  of the shift:

$$\begin{aligned} \overline{s}_k &= \beta(I - \gamma P_k)^{-1} A_k (-b_{k-1}) \\ &= \beta(I - \gamma P_k)^{-1} A_k \left[ \left( \sum_{j=k_0}^{k-1} \beta^{k-1-j} (A_{k-1} A_{k-2} \dots A_{j+1}) (-x_j) \right) + \mathcal{O}(\gamma^{k-k_0}) \right] \\ &= \sum_{j=k_0}^{k-1} \frac{\beta^{k-j}}{1-\gamma} Y_{j,k} (I - \gamma P_j) \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \end{aligned} \quad (33)$$

with

$$Y_{j,k} := (1-\gamma)(I - \gamma P_k)^{-1} A_k A_{k-1} \dots A_{j+1}.$$

Eventually, from Equations 30, 31 and 33 we get:

$$\begin{aligned} l_k \leq \sum_{j=k_0}^{k-1} \left[ \left( \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} + \frac{\beta^{k-j}}{1-\gamma} Y_{j,k} \right) (I - \gamma P_j) - \gamma^{k-j} (P_*)^{k-j} \right] \epsilon_j \\ + \mathcal{O}(\gamma^{k-k_0}). \end{aligned} \quad (34)$$

Introduce the following matrices:

$$\begin{aligned} B_{jk} &:= \frac{1-\gamma}{\gamma^{k-j}} \left[ \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} + \frac{\beta^{k-j}}{1-\gamma} Y_{j,k} \right] \\ B'_{jk} &:= \gamma B_{jk} P_j + (1-\gamma)(P_*)^{k-j}. \end{aligned}$$

**Lemma 34**  $B_{jk}$  and  $B'_{jk}$  are stochastic matrices.

**Proof** It is clear from the definition of  $X_{i,j,k}$  and  $Y_{j,k}$  that normalizing  $B_{jk}$  and  $B'_{jk}$  will give stochastic matrices. So we just need to check that their max-norm is 1.

$$\begin{aligned} \|B_{jk}\| &= \frac{(1-\gamma)}{\gamma^{k-j}} \left[ \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} + \frac{\beta^{k-j}}{1-\gamma} \right] \\ &= \frac{(1-\gamma)}{\gamma^{k-j}} \left[ \frac{\lambda\gamma}{1-\lambda\gamma} \frac{\gamma^{k-j} - \beta^{k-j}}{\gamma - \beta} + \frac{\beta^{k-j}}{1-\gamma} \right] \\ &= \frac{(1-\gamma)}{\gamma^{k-j}} \left[ \frac{\gamma^{k-j} - \beta^{k-j}}{(1-\lambda\gamma)(1-\beta)} + \frac{\beta^{k-j}}{1-\gamma} \right] \\ &= \frac{(1-\gamma)}{\gamma^{k-j}} \left[ \frac{\gamma^{k-j} - \beta^{k-j}}{1-\gamma} + \frac{\beta^{k-j}}{1-\gamma} \right] \\ &= 1. \end{aligned}$$

where we used the identities:  $\lambda\gamma = \frac{\gamma-\beta}{1-\beta}$  and  $(1-\beta)(1-\gamma\lambda) = 1-\gamma$ . Then it is also clear that  $\|B'_{jk}\| = 1$ .  $\blacksquare$

Equation 34 can be rewritten as follows:

$$\begin{aligned} l_k &\leq \sum_{j=k_0}^{k-1} \left[ \frac{\gamma^{k-j}}{1-\gamma} B_{jk}(I - \gamma P_j) - \gamma^{k-j} (P_*)^{k-j} \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \\ &= \frac{1}{1-\gamma} \sum_{j=k_0}^{k-1} \gamma^{k-j} [B_{jk} - B'_{jk}] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \end{aligned}$$

Taking the supremum limit, we see that for all  $k_0$ ,

$$\limsup_{k \rightarrow \infty} l_k \leq \frac{1}{1-\gamma} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} [B_{jk} - B'_{jk}] \epsilon_j \quad (35)$$

## C.2 When the value converges

Suppose  $\lambda$  Policy Iteration converges to some value  $v$ . Let policy  $\pi$  be the corresponding greedy policy, with stochastic matrix  $P$ . Let  $b$  be the Bellman residual of  $v$ . It is also clear that the approximation error also converges to some  $\epsilon$ . Indeed from Algorithm 3 and Equation 10, we get:

$$b = \mathcal{T}v - v = (I - \lambda\gamma P)(-\epsilon)$$

From the bound with respect to the Bellman residual (Equation 40 page 48), we can see that:

$$\begin{aligned} v_* - v^\pi &\leq [(I - \gamma P_*)^{-1} - (I - \gamma P)^{-1}] b \\ &= [(I - \gamma P)^{-1} - (I - \gamma P_*)^{-1}] (I - \lambda\gamma P) \epsilon \\ &= [(I - \gamma P)^{-1} (I - \lambda\gamma P) - (I - \gamma P_*)^{-1} (I - \lambda\gamma P)] \epsilon \\ &= [(I - \gamma P)^{-1} (I - \gamma P + \gamma P - \lambda\gamma P) - (I - \gamma P_*)^{-1} (I - \lambda\gamma P)] \epsilon \\ &= [(I + (1-\lambda)(I - \gamma P)^{-1} \gamma P + \lambda(I - \gamma P_*)^{-1} \gamma P) - (I - \gamma P_*)^{-1}] \epsilon \\ &= [((1-\lambda)(I - \gamma P)^{-1} \gamma P + \lambda(I - \gamma P_*)^{-1} \gamma P) - (I - \gamma P_*)^{-1} \gamma P_*] \epsilon \\ &= \frac{\gamma}{1-\gamma} [B_v - D] \epsilon. \end{aligned}$$

where

$$\begin{aligned} B_v &:= (1-\gamma) \left( (1-\lambda)(I - \gamma P)^{-1} P + \lambda(I - \gamma P_*)^{-1} P \right) \\ D &:= (1-\gamma) P_* (I - \gamma P_*)^{-1}. \end{aligned}$$

**Lemma 35**  $B_v$  and  $D$  are stochastic matrices.

**Proof** It is clear that  $\|D\| = 1$ . Also:

$$\begin{aligned}\|B_v\| &= (1 - \gamma) \left( 1 + \frac{(1 - \lambda)\gamma}{1 - \gamma} + \frac{\lambda\gamma}{1 - \gamma} \right) \\ &= (1 - \gamma) \left( 1 + \frac{\gamma}{1 - \gamma} \right) \\ &= 1.\end{aligned}$$

Then, the first bound of Corollary 27 follows from the application of Lemmas 23 and 25. ■

### C.3 When the policy converges

Suppose  $\lambda$  Policy Iteration converges to some policy  $\pi$ . Write  $P$  the corresponding stochastic matrix and

$$A^\pi := (1 - \lambda\gamma)P(I - \lambda\gamma P)^{-1}.$$

Then for some big enough  $k_0$ , we have:

$$l_k \leq \sum_{j=k_0}^{k-1} \left[ \frac{\gamma^{k-j}}{1 - \gamma} A^{\pi}_{jk} A^\pi (I - \gamma P) - \gamma^{k-j} (P_*)^{k-j} \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0})$$

where

$$A^{\pi}_{jk} := \frac{1 - \gamma}{\gamma^{k-j}} \left[ \frac{\lambda\gamma}{1 - \lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} (P_*)^{k-1-i} (A^\pi)^{i-j} + \beta^{k-j} (I - \gamma P)^{-1} (A^\pi)^{k-1-j} \right]$$

is a stochastic matrix (for the same reasons why  $B_{jk}$  is a stochastic matrix in Lemma 34). Noticing that

$$\begin{aligned}A^\pi (I - \gamma P) &= (1 - \lambda\gamma)P(I - \lambda\gamma P)^{-1}(I - \gamma P) \\ &= (1 - \lambda\gamma)P(I - \lambda\gamma P)^{-1}(I - \lambda\gamma P + \lambda\gamma P - \gamma P) \\ &= (1 - \lambda\gamma)P(I - (1 - \lambda)(I - \lambda\gamma P)^{-1}\gamma P) \\ &= (1 - \lambda\gamma)P - \gamma(1 - \lambda)A^\pi P\end{aligned}$$

we can deduce that

$$\begin{aligned}l_k &\leq \sum_{j=k_0}^{k-1} \left[ \frac{\gamma^{k-j}}{1 - \gamma} A^{\pi}_{jk} [(1 - \lambda\gamma)P - \gamma(1 - \lambda)A^\pi P] - \gamma^{k-j} (P_*)^{k-j} \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \\ &= \sum_{j=k_0}^{k-1} \gamma^{k-j} \left[ \frac{1 - \lambda\gamma}{1 - \gamma} A^{\pi}_{jk} P - \left[ \frac{\gamma(1 - \lambda)}{1 - \gamma} A^{\pi}_{jk} A^\pi P + (P_*)^{k-j} \right] \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \\ &= \frac{1 - \lambda\gamma}{1 - \gamma} \sum_{j=k_0}^{k-1} \gamma^{k-j} [B_{jk}^\pi - B_{jk}^{\prime\pi}] \epsilon_j + \mathcal{O}(\gamma^{k-k_0})\end{aligned}\tag{36}$$

where

$$\begin{aligned} B_{jk}^\pi &:= A^\pi_{jk} P \\ B'_{jk} &:= \frac{1-\gamma}{1-\lambda\gamma} \left[ \frac{\gamma(1-\lambda)}{1-\gamma} A^\pi_{jk} A^\pi P + (P_*)^{k-j} \right]. \end{aligned}$$

**Lemma 36**  $B_{jk}^\pi$  and  $B'_{jk}$  are stochastic matrices.

**Proof** It is clear that  $\|B_{jk}^\pi\| = 1$ . Also:

$$\begin{aligned} \|B'_{jk}\| &= \frac{1-\gamma}{1-\lambda\gamma} \left( 1 + \frac{\gamma(1-\lambda)}{1-\gamma} \right) \\ &= \frac{1-\gamma}{1-\lambda\gamma} \frac{1-\gamma + \gamma - \lambda\gamma}{1-\gamma} \\ &= 1. \end{aligned}$$

Then, the second bound of Corollary 27 follows from the application of Lemmas 23 and 25. ■

## Appendix D. Componentwise bounds with respect to the Bellman residuals

In this Appendix, we study the loss

$$l_k := v_* - v^{\pi_k}$$

with respect to the two following **Bellman residuals**:

$$b'_k := \mathcal{T}_k v_k - v_k$$

$$b_k := \mathcal{T}_{k+1} v_k - v_k = \mathcal{T} v_k - v_k$$

$b'_k$  says how much  $v_k$  differs from the value of  $\pi_k$  while  $b_k$  says how much  $v_k$  differs from the value of the policies  $\pi_{k+1}$  and  $\pi_*$ . The results stated here prove the last two inequalities of Lemma 21 page 21.

### D.1 Bounds with respect to Policy Bellman residual

Our analysis relies on the following lemma

**Lemma 37** *Suppose we have a policy  $\pi$ , a function  $v$  that is an approximation of the value  $v^\pi$  of  $\pi$  in the sense that its residual  $b' := T^\pi v - v$  is small. Taking the greedy policy  $\pi'$  with respect to  $v$  will reduce the loss as follows:*

$$v_* - v^{\pi'} \leq \gamma P_*(v_* - v^\pi) + (\gamma P_*(I - \gamma P)^{-1} - \gamma P'(I - \gamma P')^{-1}) b'$$

where  $P$  and  $P'$  are the stochastic matrices which correspond to  $\pi$  and  $\pi'$ .

**Proof** We have:

$$\begin{aligned}
 v_* - v^{\pi'} &= \mathcal{T}_* v_* - \mathcal{T}^{\pi'} v^{\pi'} \\
 &= \mathcal{T}_* v_* - \mathcal{T}_* v^\pi + \mathcal{T}_* v^\pi - \mathcal{T}_* v + \mathcal{T}_* v - \mathcal{T}^{\pi'} v + \mathcal{T}^{\pi'} v - \mathcal{T}^{\pi'} v^{\pi'} \\
 &\leq \gamma P_*(v_* - v^\pi) + \gamma P_*(v^\pi - v) + \gamma P'(v - v^{\pi'})
 \end{aligned} \tag{37}$$

where we used the fact that  $\mathcal{T}_* v \leq \mathcal{T}^{\pi'} v$ . One can see that:

$$\begin{aligned}
 v^\pi - v &= \mathcal{T}^\pi v^\pi - v \\
 &= \mathcal{T}^\pi v^\pi - \mathcal{T}^\pi v + \mathcal{T}^\pi v - v \\
 &= \gamma P(v^\pi - v) + b' \\
 &= (I - \gamma P)^{-1} b'
 \end{aligned} \tag{38}$$

and that

$$\begin{aligned}
 v - v^{\pi'} &= v - \mathcal{T}^{\pi'} v^{\pi'} \\
 &= v - \mathcal{T}^\pi v + \mathcal{T}^\pi v - \mathcal{T}^{\pi'} v + \mathcal{T}^{\pi'} v - \mathcal{T}^{\pi'} v^{\pi'} \\
 &\leq -b' + \gamma P'(v - v^{\pi'}) \\
 &\leq (I - \gamma P')^{-1}(-b').
 \end{aligned} \tag{39}$$

where we used the fact that  $\mathcal{T}^\pi v \leq \mathcal{T}^{\pi'} v$ . We get the result by putting back Equations 38 and 39 into Equation 37.  $\blacksquare$

To derive a bound for  $\lambda$  Policy Iteration, we simply apply the above lemma to  $\pi = \pi_k$ ,  $v = v_k$  and  $\pi' = \pi_{k+1}$ . We thus get:

$$l_{k+1} \leq \gamma P_* l_k + (\gamma P_*(I - \gamma P_k)^{-1} - \gamma P_{k+1}(I - \gamma P_{k+1})^{-1}) b'_k$$

Introduce the following stochastic matrices:

$$\begin{aligned}
 C_k &:= (1 - \gamma)^2 (I - \gamma P_*)^{-1} (P_*(I - \gamma P_k)^{-1}) \\
 C'_k &:= (1 - \gamma)^2 (I - \gamma P_*)^{-1} (P_{k+1}(I - \gamma P_{k+1})^{-1})
 \end{aligned}$$

This leads to the following componentwise bound:

$$\limsup_{k \rightarrow \infty} l_k \leq \frac{\gamma}{(1 - \gamma)^2} \limsup_{k \rightarrow \infty} [C_k - C'_k] b'_k$$

## D.2 Bounds with respect to Bellman residual

We rely on the following lemma (which is for instance proved by Munos in Munos (2007))

**Lemma 38** *Suppose we have a function  $v$ . Let  $\pi$  be the greedy policy with respect to  $v$ . Then*

$$v_* - v^\pi \leq \gamma [P_*(I - \gamma P_*)^{-1} - P^\pi(I - \gamma P^\pi)^{-1}] (\mathcal{T}^\pi v - v)$$



We provide a proof for the sake of completeness:

**Proof** Using the fact that  $\mathcal{T}_*v \leq \mathcal{T}^\pi v$ , we see that

$$\begin{aligned}
v_* - v^\pi &= \mathcal{T}_*v_* - \mathcal{T}^\pi v^\pi \\
&= \mathcal{T}_*v_* - \mathcal{T}_*v + \mathcal{T}_*v - \mathcal{T}^\pi v + \mathcal{T}^\pi v - \mathcal{T}^\pi v^\pi \\
&\leq \mathcal{T}_*v_* - \mathcal{T}_*v + \mathcal{T}^\pi v - \mathcal{T}^\pi v^\pi \\
&= \gamma P_*(v_* - v) + \gamma P^\pi(v - v^\pi) \\
&= \gamma P_*(v_* - v^\pi) + \gamma P_*(v^\pi - v) + \gamma P^\pi(v - v^\pi) \\
&\leq (I - \gamma P_*)^{-1}(\gamma P_* - \gamma P^\pi)(v^\pi - v).
\end{aligned}$$

Using Equation 38 we see that:

$$v^\pi - v = (I - \gamma P^\pi)^{-1}(\mathcal{T}^\pi v - v).$$

Thus

$$\begin{aligned}
v_* - v^\pi &\leq (I - \gamma P_*)^{-1}(\gamma P_* - \gamma P^\pi)(I - \gamma P^\pi)^{-1}(\mathcal{T}^\pi v - v) \\
&= (I - \gamma P_*)^{-1}(\gamma P_* - I + I - \gamma P^\pi)(I - \gamma P^\pi)^{-1}(\mathcal{T}^\pi v - v) \\
&= [(I - \gamma P_*)^{-1} - (I - \gamma P^\pi)^{-1}](\mathcal{T}^\pi v - v) \\
&= \gamma [P_*(I - \gamma P_*)^{-1} - P^\pi(I - \gamma P^\pi)^{-1}](\mathcal{T}^\pi v - v). \quad \blacksquare
\end{aligned}$$

To derive a bound for  $\lambda$  Policy Iteration, we simply apply the above lemma to  $v = v_{k-1}$  and  $\pi = \pi_k$ . We thus get:

$$l_k \leq \frac{\gamma}{1 - \gamma} [D - D'_k] b_{k-1} \quad (40)$$

where

$$\begin{aligned}
D &:= (1 - \gamma)P_*(I - \gamma P_*)^{-1} \\
D'_k &:= (1 - \gamma)P_k(I - \gamma P_k)^{-1}
\end{aligned}$$

are stochastic matrices.

## Appendix E. From componentwise bounds to span seminorm bounds

This Appendix contains the proofs of Lemmas 18 and 23 that enable us to derive span seminorm performance bounds from the componentwise analysis developed in the previous Appendices. It is easy to see that Lemma 18 is a special case of Lemma 23, so we will only prove the latter.

Write  $a_{kj} := \arg \min_a \|y_j - ae\|_{p, \mu_{kj}}$ . As  $X_{jk}$  and  $X'_{jk}$  are stochastic matrices,  $X_k e = X'_k e = e$  and we can write that:

$$\limsup_{k \rightarrow \infty} |x_k| \leq K \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \xi_{k-j} (X_{kj} - X'_{kj})(y_j - a_{kj}e).$$

By taking the absolute value we get

$$\limsup_{k \rightarrow \infty} |x_k| \leq K \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \xi_{k-j} (X_{kj} + X'_{kj}) |y_j - a_{kj}e|.$$

It can then be seen that

$$\begin{aligned} & \limsup_{k \rightarrow \infty} \left( \|x_k\|_{p,\mu} \right)^p \\ = & K^p \limsup_{k \rightarrow \infty} \mu^T (|x_k|)^p \\ \leq & K^p \limsup_{k \rightarrow \infty} \mu^T \left[ \sum_{j=k_0}^{k-1} \xi_{k-j} (X_{kj} + X'_{kj}) (|y_j - a_{kj}e|) \right]^p \\ = & K^p \limsup_{k \rightarrow \infty} \mu^T \left[ \frac{\left( \sum_{j=k_0}^{k-1} \xi_{k-j} \frac{1}{2} (X_{kj} + X'_{kj}) 2 (|y_j - a_{kj}e|) \right)}{\sum_{j=k_0}^{k-1} \xi_{k-j}} \right]^p \left( \sum_{j=k_0}^{k-1} \xi_{k-j} \right)^p \end{aligned}$$

By using Jensen's inequality (with the convex function  $x \mapsto x^p$ ), we get:

$$\begin{aligned} & \limsup_{k \rightarrow \infty} \left( \|x_k\|_{p,\mu} \right)^p \\ \leq & K^p \limsup_{k \rightarrow \infty} \mu^T \frac{\sum_{j=k_0}^{k-1} \xi_{k-j} \frac{1}{2} (X_{kj} + X'_{kj}) [2|y_j - a_{kj}e|]^p}{\sum_{j=k_0}^{k-1} \xi_{k-j}} \left( \sum_{j'=k_0}^{k-1} \xi_{k-j'} \right)^p \\ = & K^p \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \xi_{k-j} \mu_{k_j}^T [2|y_j - a_{kj}e|]^p \left( \sum_{j'=k_0}^{k-1} \xi_{k-j'} \right)^{p-1} \\ \leq & K^p \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \xi_{k-j} \left[ 2 \|y_j - a_{kj}e\|_{p,\mu_{k_j}} \right]^p K'^{p-1} \\ = & K^p K'^{p-1} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \xi_{k-j} \left[ \text{span}_{p,\mu_{k_j}} [y_j] \right]^p \\ \leq & K^p K'^{p-1} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \xi_{k-j} \left[ \sup_{k' \geq j' \geq k_0} \text{span}_{p,\mu_{k'j'}} [y_{j'}] \right]^p \\ = & K^p K'^{p-1} K' \left[ \sup_{k' \geq j' \geq k_0} \text{span}_{p,\mu_{k'j'}} [y_{j'}] \right]^p \\ = & K^p K'^p \left[ \sup_{k' \geq j' \geq k_0} \text{span}_{p,\mu_{k'j'}} [y_{j'}] \right]^p \end{aligned}$$

where we used  $\sum_{j=k_0}^{k-1} \xi_{k-j} = K'$ .

As this is true for all  $k_0$ , and as  $k_0 \mapsto \sup_{k' \geq j' \geq k_0} \text{span}_{p,\mu_{k'j'}} [y_{j'}]$  is non-increasing, the result follows.

## Appendix F. Analysis of the Undiscounted Case

This last appendix contains the Proofs of Lemma 28 and Proposition 30 that provide the analysis of an undiscounted problem like Tetris.

### F.1 Proof of Lemma 28

First of all, the relation expressed in Equation 21 between the loss and the stochastic matrices, we restate here for clarity:

$$\forall k_0, \quad \limsup_{k \rightarrow \infty} v_* - v^{\pi_k} \leq \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \delta_{k-j} [G_{jk} - G'_{jk}] \epsilon_j,$$

is obtained by simply rewriting the first inequality of Lemma 21 with  $\gamma = 1$  and  $\beta = 1$  (note in particular that the terms  $\delta_{k-j}$  collapse through the definition of  $G_{jk}$  and  $G'_{jk}$ ).

To complete the proof of the Lemma, we need to show that the matrices  $G_{jk}$  and  $G'_{jk}$  are substochastic matrices. By construction, these matrices are sum of non-negative matrices so we only need to show that their max-norm is smaller than or equal to 1.

For all  $n$ , write  $\mathcal{M}_n$  the set of matrices that is defined as follows:

- For all sets of  $n$  policies  $(\pi_1, \pi_2, \dots, \pi_n)$ ,  $P_{\pi_1} P_{\pi_2} \dots P_{\pi_n} \in \mathcal{M}_n$ ;
- For all  $\eta \in (0, 1)$ , and  $(P, Q) \in \mathcal{M}_n \times \mathcal{M}_n$ ,  $\eta P + (1 - \eta)Q \in \mathcal{M}_n$ .

The motivation for introducing this set is that we have the following properties: For all  $n$ ,  $P \in \mathcal{M}_n$  is a substochastic matrix such that  $\|P\|_\infty \leq \alpha^{\lfloor \frac{n}{n_0} \rfloor}$ . We will use the somewhat abusive notation  $\Pi_n$  for denoting any element of  $\mathcal{M}_n$ . For instance, for some matrix  $P$ , writing  $P = a\Pi_i + b\Pi_j\Pi_k$  should be read as follows: there exists  $P_1 \in \mathcal{M}_i$ ,  $P_2 \in \mathcal{M}_j$  and  $P_3 \in \mathcal{M}_k$  such that  $P = aP_1 + bP_2P_3$ .

Recall the definition of the substochastic matrix

$$A_k = (1 - \lambda)(I - \lambda P_k)P_k = (1 - \lambda) \sum_{i=0}^{\infty} \lambda^i \Pi_{i+1}.$$

Let  $i \leq j < k$ . It can be seen that

$$\begin{aligned} (P_*)^{k-1-i} A_{i+1} A_i \dots A_{j+1} &= \underbrace{\Pi_{k-1-i} \left( (1 - \lambda) \sum_{i=0}^{\infty} \lambda^i \Pi_{i+1} \right) \dots \left( (1 - \lambda) \sum_{i=0}^{\infty} \lambda^i \Pi_{i+1} \right)}_{i-j+1 \text{ terms}} \\ &= \underbrace{\Pi_{k-j} \left( (1 - \lambda) \sum_{i=0}^{\infty} \lambda^i \Pi_i \right) \dots \left( (1 - \lambda) \sum_{i=0}^{\infty} \lambda^i \Pi_i \right)}_{i-j+1 \text{ terms}}. \end{aligned} \quad (41)$$

Now, observe that

$$\begin{aligned}
 \left\| \sum_{i=0}^{\infty} \lambda^i \Pi_i \right\|_{\infty} &\leq \sum_{i=0}^{\infty} \lambda^i \|\Pi_i\|_{\infty} \\
 &\leq \sum_{i=0}^{\infty} \lambda^i \alpha^{\lfloor \frac{i}{n_0} \rfloor} \\
 &= \sum_{j=0}^{\infty} \sum_{i=0}^{n_0-1} \lambda^{j n_0 + i} \alpha^j \\
 &= \sum_{j=0}^{\infty} (\lambda^{n_0} \alpha)^j \sum_{i=0}^{n_0-1} \lambda^i \\
 &= \frac{1 - \lambda^{n_0}}{(1 - \lambda^{n_0} \alpha)(1 - \lambda)}.
 \end{aligned} \tag{42}$$

As a consequence, writing  $\eta := \frac{1 - \lambda^{n_0}}{1 - \lambda^{n_0} \alpha}$ , we see from Equations 41 that

$$\left\| (P_*)^{k-1-i} A_{i+1} A_i \dots A_{j+1} \right\|_{\infty} \leq \alpha^{\lfloor \frac{k-j}{n_0} \rfloor} \eta^{i-j+1}.$$

Similarly, by using Equation 42 and noticing that  $\frac{1 - \lambda^{n_0}}{1 - \lambda} \xrightarrow{\lambda \rightarrow 1} n_0$ , it can be seen that

$$\left\| (I - P_k)^{-1} A_k A_{k-1} \dots A_{j+1} \right\|_{\infty} \leq \frac{n_0}{1 - \alpha} \alpha^{\lfloor \frac{k-j}{n_0} \rfloor} \eta^{k-j}.$$

We are ready to bound the norm of the matrix  $G_{jk}$ :

$$\begin{aligned}
 \|G_{jk}\|_{\infty} &\leq \frac{\alpha^{\lfloor \frac{k-j}{n_0} \rfloor}}{\delta_{k-j}} \left[ \frac{\lambda}{1 - \lambda} \sum_{i=j}^{k-1} \eta^{i-j+1} + \frac{n_0 \eta^{k-j}}{1 - \alpha} \right] \\
 &= \frac{\alpha^{\lfloor \frac{k-j}{n_0} \rfloor}}{\delta_{k-j}} \left[ \left( \frac{\lambda}{1 - \lambda} \right) \eta \left( \frac{1 - \eta^{k-j}}{1 - \eta} \right) + \frac{n_0 \eta^{k-j}}{1 - \alpha} \right] \\
 &= \frac{\alpha^{\lfloor \frac{k-j}{n_0} \rfloor}}{\delta_{k-j}} \left[ \left( \frac{\lambda}{1 - \lambda} \right) \left( \frac{1 - \lambda^{n_0}}{1 - \lambda^{n_0} \alpha} \right) \left( \frac{1 - \eta^{k-j}}{1 - \eta} \right) + \frac{n_0 \eta^{k-j}}{1 - \alpha} \right] \\
 &= \frac{\alpha^{\lfloor \frac{k-j}{n_0} \rfloor}}{\delta_{k-j}} \left[ \left( \frac{1 - \lambda^{n_0}}{1 - \lambda} \right) \left( \frac{\lambda}{1 - \lambda^{n_0} \alpha} \right) \left( \frac{1 - \eta^{k-j}}{1 - \eta} \right) + \frac{n_0 \eta^{k-j}}{1 - \alpha} \right] \\
 &= 1.
 \end{aligned}$$

where we used the definition of  $\eta$ . Therefore  $G_{jk}$  is a substochastic matrix. It trivially follows that  $G'_{jk}$  is also a substochastic matrix.

## F.2 Proof of Proposition 30

In order to prove the  $L_p$  norm bound of Proposition 30, we will rely on the following variation of Lemma 23.

**Lemma 39** *If  $x_k$  and  $y_k$  are sequences of vectors and  $X_{jk}, X'_{jk}$  sequences of substochastic matrices satisfying*

$$\forall k_0, \quad \limsup_{k \rightarrow \infty} \|x_k\| \leq K \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \xi_{k-j} (X_{kj} - X'_{kj}) y_j,$$

where  $(\xi_i)_{i \geq 1}$  is a sequence of non-negative weights satisfying:

$$\sum_{i=1}^{\infty} \xi_i = K' < \infty,$$

then, for all distribution  $\mu$ ,

$$\mu_{kj} := \frac{1}{2} (X_{kj} + X'_{kj})^T \mu$$

is a non-negative vector and  $\tilde{\mu}_{kj} := \frac{\mu_{kj}}{\|\mu_{kj}\|_1}$  is a distribution, and

$$\limsup_{k \rightarrow \infty} \|x_k\|_{p,\mu} \leq KK' \lim_{k_0 \rightarrow \infty} \left[ \sup_{k \geq j \geq k_0} \|y_j\|_{p,\tilde{\mu}_{kj}} \right].$$

**Proof** The proof follows the lines of that of Lemma 23 in Appendix E. The differences are as follows:

- Since  $X_{jk}$  and  $X'_{jk}$  are substochastic matrices (and not stochastic matrices), we have in general  $X_{jk}e \neq X'_{jk}e$  and must take  $a_{kj} = 0$ , which in turn gives an  $L_p$  norm bound instead of the  $L_p$  span seminorm bound.
- To express the bound in terms of the distributions  $\tilde{\mu}_{kj}$ , we use the fact that  $\mu_{kj} \leq \tilde{\mu}_{kj}$  which derives from  $\|\mu_{kj}\|_1$  since  $X_{jk}$  and  $X'_{jk}$  are substochastic matrices. ■

Proposition 30 is obtained by applying this Lemma and an analogue of Lemma 25 for  $L_p$  norm on the componentwise bound (Lemma 28 - see previous subsection). The only remaining thing that needs to be checked is that  $\sum_{i=1}^{\infty} \delta_i$  has the right value. This is what we do now.

Similarly to Equation 42, one can see that:

$$\sum_{i=0}^{\infty} \alpha^{\lfloor \frac{i}{n_0} \rfloor} \eta^i = \frac{1 - \eta^{n_0}}{(1 - \eta^{n_0} \alpha)(1 - \eta)}$$

and

$$\sum_{i=0}^{\infty} \alpha^{\lfloor \frac{i}{n_0} \rfloor} (1 - \eta^i) = \frac{n_0}{1 - \alpha} - \frac{1 - \eta^{n_0}}{(1 - \eta^{n_0} \alpha)(1 - \eta)}.$$

As a consequence:

$$\begin{aligned}
 \sum_{i=0}^{\infty} \delta_i &= \sum_{i=0}^{\infty} \alpha^{\lfloor \frac{i}{n_0} \rfloor} \left( \frac{1 - \lambda^{n_0}}{1 - \lambda} \right) \left( \frac{\lambda}{1 - \lambda^{n_0} \alpha} \right) \left( \frac{1 - \eta^i}{1 - \eta} \right) + \frac{n_0 \eta^i}{1 - \alpha} \\
 &= \left( \frac{1 - \lambda^{n_0}}{1 - \lambda} \right) \left( \frac{\lambda}{1 - \lambda^{n_0} \alpha} \right) \left( \frac{\sum_{i=0}^{\infty} \alpha^{\lfloor \frac{i}{n_0} \rfloor} (1 - \eta^i)}{1 - \eta} \right) + \frac{n_0 \sum_{i=0}^{\infty} \alpha^{\lfloor \frac{i}{n_0} \rfloor} \eta^i}{1 - \alpha} \\
 &= \left( \frac{1 - \lambda^{n_0}}{1 - \lambda} \right) \left( \frac{\lambda}{1 - \lambda^{n_0} \alpha} \right) \left( \frac{1}{1 - \eta} \right) \left( \frac{n_0}{1 - \alpha} - \frac{1 - \eta^{n_0}}{(1 - \eta^{n_0} \alpha)(1 - \eta)} \right) \\
 &\quad + \left( \frac{n_0}{1 - \alpha} \right) \left( \frac{1 - \eta^{n_0}}{(1 - \eta^{n_0} \alpha)(1 - \eta)} \right) \\
 &= \lambda f(\lambda) \frac{1}{1 - \eta} (f(1) - f(\eta)) + f(1) f(\eta) \tag{43}
 \end{aligned}$$

with for all  $x$ ,  $f(x) := \frac{(1-x^{n_0})}{(1-x)(1-x^{n_0}\alpha)}$  and  $f(1) = \frac{n_0}{1-\alpha}$  by continuity. Now, we can conclude by noticing that

$$\sum_{i=1}^{\infty} \delta_i = \sum_{i=0}^{\infty} \delta_i - \delta_0$$

and  $\delta_0 = \frac{n_0}{1-\alpha} = f(1)$ .