



HAL
open science

Performance bounds for Lambda Policy Iteration

Bruno Scherrer

► **To cite this version:**

Bruno Scherrer. Performance bounds for Lambda Policy Iteration. [Research Report] 2007, pp.29.
inria-00185271v1

HAL Id: inria-00185271

<https://inria.hal.science/inria-00185271v1>

Submitted on 5 Nov 2007 (v1), 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 λ Policy Iteration

Bruno Scherrer

5 novembre 2007

1 Abstract

We consider the discrete-time infinite-horizon discounted stationary optimal control problem formalized by Markov Decision Processes. We study λ Policy Iteration, a family of algorithms parameterized by λ , originally introduced by Ioffe and Bertsekas [1]. λ Policy Iteration generalizes the standard algorithms Value Iteration and Policy Iteration, and has some connections with TD(λ) introduced by Sutton & Barto [7]. We deepen the original theory developed by Ioffe and Bertsekas [1] by providing convergence rate bounds which generalize standard bounds for Value Iteration described for instance by Puterman [6]. We also develop the theory of this algorithm when it is used in an approximate form. Doing so, we extend and unify the separate analyses developed by Munos for Approximate Value Iteration [5] and Approximate Policy Iteration [4].

2 Definition of norms and seminorms

The analysis we will describe in this article relies on several norms and seminorms which we define here.

Let X be a finite space. In this section, u denotes 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. μ denotes a distribution on X . We consider the weighted L_p norm :

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

We will write $\|\cdot\|_p$ the unweighted L_p norm (i.e. with uniform distribution μ). $\|\cdot\|_\infty$ 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 [6]) :

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

It can be seen that

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

It is thus natural to generalize the span seminorm definition 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.

Several bounds we give 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 bounds with the (more standard) L_p and max norms from the a bound with the span seminorm. However, the reverse is not true : as it is not a norm (as it can in particular be 0 for non zero functions), the span seminorm can be arbitrarily smaller than these standard norms (and this constitutes our motivation for using it).

3 Markov Decision Processes

3.1 The problem

This paper is about the discrete-time infinite-horizon discounted stationary optimal control problem, which we describe now. 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 , 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** [6].

We are interested in stationary deterministic policies, that is functions $\pi : X \rightarrow A$ which map states into controls¹. Writing i_k the state at time k , the **value of policy** π at state i is defined as the total expected return while following a policy π from i , that is

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

where E_π denotes with expectation conditional on the fact that the actions are selected with the policy π . The **optimal value** starting from state i is defined as

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

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

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

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

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

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

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

v_* is a fixed point of the nonlinear backup operator $\mathcal{T}v := \max_\pi \mathcal{T}^\pi v$. Once the optimal value v_* is computed, deriving an optimal policy is straightforward. For any value vector v , we call a **greedy policy with respect to the value** v a policy π 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 notation² $\text{greedy}(v)$ any policy that is greedy with respect to v . The notions of optimal value function and greedy policies are fundamental to optimal

¹For the criterion to be defined soon that we wish to optimize, it can indeed be shown that there exists at least one stationary deterministic policy which is optimal [6].

²There might be several policies that are greedy with respect to a value v .

control because of the following property : any policy π_* that is greedy with respect to the optimal value is an **optimal policy** and its value v^{π_*} equals v_* .

It is well-known that the backup operators \mathcal{T}^π and \mathcal{T} are γ -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}^π . Being a γ -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 implies that the corresponding fixed point exists and is unique and that for any initial vector v_0 ,

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

3.2 Value Iteration

The Value Iteration algorithms for computing the value of a policy π and the value of the optimal policy π_* rely on equation 4. Algorithm 1 provides a description of Value Iteration for computing an optimal policy (replace \mathcal{T} by \mathcal{T}^π in it and one gets Value Iteration for computing the value of policy π). In this description,

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)

we have introduced a term ϵ_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^{π_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 π_* (see for instance [6]). Let π_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 $\|v_{k_0+1} - v_{k_0}\|_\infty < \frac{1-\gamma}{2\gamma}\epsilon$, then $\|v_* - v^{\pi_{k_0+1}}\|_\infty < \epsilon$.

It is somewhat less known that the Bellman operators \mathcal{T} and \mathcal{T}^π are also contraction mapping with respect to the span seminorm [6]. 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 [6] :

Theorem 1 (Stopping Condition for Exact Value Iteration [6])

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

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

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

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

Properties of Approximate Value Iteration 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 [2] : If the approximation errors are uniformly bounded $\|\epsilon_k\|_\infty \leq \epsilon$, then the difference between the asymptotic performance of policies π_{k+1} greedy with respect to v_k and the optimal policy satisfies

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

Munos has recently argued [5] that, since most supervised learning algorithms (such as least square regression) that are used in practice for approximating each iterate of Value Iteration minimize an empirical approximation in some L_p norm, it would be more interesting to have an extension of the above result where the approximation error ϵ is expressed in terms of the L_p norm. Munos actually showed how to do this in [5]. The idea is to analyze the *componentwise* asymptotic behaviour of Approximate Value Iteration, from which it is rather easy to derive L_p analysis for any p .

Lemma 1 (Componentwise Asymptotic Performance of Approximate Value Iteration [5])

Write $P_k = P^{\pi_k}$ the stochastic matrix corresponding to the policy π_k which is greedy with respect to v_{k-1} , P_* the stochastic matrix corresponding to the (unknown) optimal policy π_* , and denote $|\epsilon_k|$ the vector of absolute values of the ϵ_k . Then :

$$\limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty \leq \limsup_{k \rightarrow \infty} (I - \gamma P_k)^{-1} \sum_{j=0}^{k-1} \gamma^{k-j} [P_k P_{k-1} \dots P_{j+1} - (P_*)^{k-j}] |\epsilon_j|.$$

Munos introduces some **concentration coefficient** [4, 5] : Assume there exists a distribution ν and a real number $C(\nu)$ such that

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

For instance, if one chooses the uniform law ν , then there always exists such a $C(\nu) \in (1, N)$ where N is the size of the state space. See [4, 5] for more discussion on this coefficient. From the above Lemma it is possible to show that :

Theorem 2 (Asymptotic Performance of Approximate Value Iteration in L_p norm [5])

If the approximation errors are uniformly bounded $\|\epsilon_k\|_{p,\nu} \leq \epsilon$, then the difference between the asymptotic performance of policies π_{k+1} greedy with respect to v_k and the optimal policy satisfies

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

The main difference with the previous bound by Bertsekas and Tsitsiklis (equation 5) is that the approximation error ϵ is controlled by than L_p weighted norm. As $(C(\nu))^{1/p} \xrightarrow{p \rightarrow \infty} 1$ this results is strictly better.

There is in general no guarantee that AVI converges. AVI may be shown to converge in some specific cases such as state aggregation and when using the so-called averagers [3]. Also, convergence may just occur experimentally. Suppose (v_k) tends to some v . Write π the corresponding greedy policy. Note also that (ϵ_k) tends to $v - \mathcal{T}v$, which is known as the **Bellman residual**. The above bounds apply, but in this specific case, they can be improved by a factor $\frac{1}{1-\gamma}$. It is indeed known (e.g. [8]) that

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

and, with the same notations as above, Munos derived the analogous better L_p bound [5] :

Corollary 1 (Performance of Approximate Value Iteration in case of convergence [5])

Suppose (v_k) tends to some v . Write π the corresponding greedy policy. Then

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

Eventually, let us mention that in [5], Munos also shows some *finer* performance bounds (in L_p weighted norm) using some *finer* concentration coefficients. As we won't discuss this in this paper, we recommend the interested reader to go through [5].

3.3 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 property : if π is some policy,

Algorithm 2 Policy Iteration

Input: An MDP, an initial policy π_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 π_k

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

$k \leftarrow k + 1$

until some stopping criterion

Return π_k

then any policy π' that is greedy with respect to the value of π , i.e. any π' satisfying $\pi' = \text{greedy}(v^\pi)$, is better than π in the sense that $v^{\pi'} \geq v^\pi$. Policy Iteration iterates this process for improving a policy in order to generate a sequence of policies with increasing values. It is described in Algorithm 2. We use the value of a policy given by equation 3. As for Value Iteration, our description includes a potential error ϵ_k term each time the value of a policy is estimated.

Properties of Exact Policy Iteration When the the state space and the control spaces are finite, it is well-known that Exact Policy Iteration converges to an optimal policy π_* in a finite number of iterations. It

is known that the rate of convergence is at least linear [6]. If the function $v \mapsto P^{\text{greedy}(v)}$ is Lipschitz, then it can be shown that Policy Iteration has a quadratic convergence [6]. However, we did not find any stopping condition in the literature that is similar to the one of Theorem 1.

Properties of Approximate Policy Iteration 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 analog of equation 5 which is also due to Bertsekas and Tsitsiklis [2] : If the approximation errors are uniformly bounded $\|\epsilon_k\|_\infty \leq \epsilon$, then the difference between the asymptotic performance of policies π_{k+1} greedy with respect to v_k and the optimal policy is

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

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 π_k ; this is interesting when the evaluation step of Approximate Policy Iteration involves the minimization of this Bellman residual. 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 π_{k+1} is greedy with respect to v_k). To avoid any 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 :

Lemma 2 (Componentwise Asymptotic Performance of Approximate Policy Iteration [4])

Write $P_k = P^{\pi_k}$ the stochastic matrix corresponding to the policy π_k which is greedy with respect to v_{k-1} , P_* the stochastic matrix corresponding to the (unknown) optimal policy π_* , and denote $|\epsilon_k|$ the vector of absolute values of the ϵ_k . Then :

$$\begin{aligned} \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty &\leq (I - \gamma P_*)^{-1} \limsup_{k \rightarrow \infty} [\gamma P_{k+1} (I - \gamma P_{k+1})^{-1} (I - \gamma P_k) - \gamma P_*] \epsilon_k \\ \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty &\leq (I - \gamma P_*)^{-1} \limsup_{k \rightarrow \infty} [\gamma P_{k+1} (I - \gamma P_{k+1})^{-1} - \gamma P_* (I - \gamma P_k)] (v_k - \mathcal{T}^{\pi_k} v_k). \end{aligned}$$

Using the concentration coefficient $C(\nu)$ introduced in the previous section (equation 6), it is possible to show³ the following L_p bounds :

Theorem 3 (Asymptotic Performance of Approximate Policy Iteration in L_p norm)

The difference between the asymptotic performance of policies π_{k+1} greedy with respect to v_k and the optimal policy satisfies

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

Note that the bound with respect to the approximation error ϵ is better than Bertsekas and Tsitsiklis's (equation 7) Compared to the analog result for Approximate Value Iteration (Theorem 2) 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 (from [4], Remark 4 page 7) that :

³This result is not stated by Munos in [4] but using techniques of another of his articles [5], it is straightforward to derive from Lemma 2. The current paper will anyway generalize this result (in Theorem 9 page 19).

Corollary 2 (Performance of Approximate Policy Iteration in case of convergence)

If the sequence of policies (π_k) converges to some π , then

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

3.4 λ Policy Iteration

Value Iteration and Policy Iteration are often considered as unrelated algorithms. Though all the results we have emphasized so far are strongly related (and even sometimes identical), they were proven independently for each algorithm. In this section, we describe a family of algorithms called “ λ Policy Iteration” (originally introduced in [1]) parameterized by a coefficient $\lambda \in (0, 1)$, that generalizes them both. When $\lambda = 0$, λ Policy Iteration will reduce to Value Iteration while it will reduce to Policy Iteration when $\lambda = 1$. We will also briefly discuss the fact that λ Policy Iteration draws some connections with Temporal Difference algorithms which one finds in the Reinforcement Learning literature [7].

We begin by giving some intuition about how one can make a connection between Value Iteration and Policy Iteration. For the moment let us forget about the error term ϵ_k . Though Value Iteration seems to build a sequence of value functions and Policy Iteration to build a sequence of policies, both algorithms can in fact 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 π_{k+1} is the fixed point of $\mathcal{T}^{\pi_{k+1}}$ — it can be argued that every iteration of Policy Iteration does the following :

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

Thanks to this little rewriting, both algorithms now look close to each other. Both can be seen as having an estimate v_k of the value of policy π_k , from which they deduce a potentially better policy π_{k+1} . The corresponding value $v^{\pi_{k+1}}$ of this better policy may be regarded as a target which is going to be 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 π_{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). We made this rewriting because it leads to a natural introduction of λ Policy Iteration : namely λ Policy Iteration is doing a λ -adjustable step towards the value of π_{k+1} :

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

More formally, the idea of λ Policy Iteration (see the above left hand side) consists in doing a 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 λ Policy Iteration generalizes Value Iteration (when $\lambda = 0$) and Policy Iteration (when $\lambda = 1$). The actual equivalence between the left and the right can be proven as

follows :

$$\begin{aligned}
(1 - \lambda) \sum_{j=0}^{\infty} \lambda^j (\mathcal{T}^{\pi_{k+1}})^{j+1} v_k &= (1 - \lambda) \sum_{j=0}^{\infty} \lambda^j \left\{ \left[\sum_{l=0}^j (\gamma P^{\pi_{k+1}})^l \right] r^{\pi_{k+1}} + (\gamma P^{\pi_{k+1}})^{j+1} v_k \right\} \\
&= \sum_{j=0}^{\infty} \sum_{l=0}^j (1 - \lambda) \lambda^j (\gamma P^{\pi_{k+1}})^l r^{\pi_{k+1}} + (1 - \lambda) \sum_{j=0}^{\infty} \lambda^j (\gamma P^{\pi_{k+1}})^{j+1} v_k \\
&= \sum_{l=0}^{\infty} \sum_{j=l}^{\infty} (1 - \lambda) \lambda^j (\gamma P^{\pi_{k+1}})^l r^{\pi_{k+1}} + (1 - \lambda) \sum_{j=0}^{\infty} \lambda^j (\gamma P^{\pi_{k+1}})^{j+1} v_k \\
&= \sum_{l=0}^{\infty} \left[\sum_{j=l}^{\infty} \left(\lambda^j (\gamma P^{\pi_{k+1}})^l r^{\pi_{k+1}} - \lambda^{j+1} (\gamma P^{\pi_{k+1}})^l r^{\pi_{k+1}} \right) \right] + (1 - \lambda) \sum_{j=0}^{\infty} \lambda^j (\gamma P^{\pi_{k+1}})^{j+1} v_k \\
&= \sum_{l=0}^{\infty} \lambda^l (\gamma P^{\pi_{k+1}})^l r^{\pi_{k+1}} + (1 - \lambda) \sum_{j=0}^{\infty} \lambda^j (\gamma P^{\pi_{k+1}})^{j+1} v_k \\
&= \sum_{l=0}^{\infty} (\lambda \gamma P^{\pi_{k+1}})^l (r^{\pi_{k+1}} + (1 - \lambda) \gamma P^{\pi_{k+1}} v_k) \\
&= (I - \lambda \gamma P^{\pi_{k+1}})^{-1} (r^{\pi_{k+1}} + (1 - \lambda) \gamma P^{\pi_{k+1}} v_k).
\end{aligned}$$

In order to describe λ Policy Iteration, it is useful to introduce a new operator. For any value v and any policy π , define (the following four formulations are equivalent up to some little linear algebra manipulations) :

$$\mathcal{T}_{\lambda}^{\pi} v := v + (I - \lambda \gamma P^{\pi})^{-1} (\mathcal{T}^{\pi} v - v) \quad (8)$$

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

$$= (I - \lambda \gamma P^{\pi})^{-1} (r^{\pi} + (1 - \lambda) \gamma P^{\pi} v) \quad (9)$$

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

λ Policy Iteration is formally described in Algorithm 3. Once again, our description includes a potential

Algorithm 3 λ 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 π_{k+1}

$k \leftarrow k + 1$

until some convergence criterion

Return $\text{greedy}(v_k)$

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 Iteration⁴ (Algorithm 2) when $\lambda = 1$.

⁴Policy Iteration starts with an initial policy while λ 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 .

Relation with Reinforcement Learning The definition of the operator \mathcal{T}_λ^π given by equation 9 is the form we have used for the introduction of λ Policy Iteration as an intermediate between Value Iteration and Policy Iteration. The equivalent form given by equation 8 can be used to make a connection with the TD(λ) algorithms⁵ that one finds in the Reinforcement Learning literature, of which a reference is the book by Sutton & Barto [7]. Indeed, equation 9 can be seen as an incremental additive procedure :

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

where the ‘‘increment’’ is $\delta_k := (I - \lambda\gamma P^{\pi_{k+1}})^{-1}(r^{\pi_{k+1}} + (1 - \lambda)\gamma P^{\pi_{k+1}}v_k)$. It can be shown (see [1] for a proof or simply look at the equivalence between equations 2 and 3 for the intuition) that the vector δ_k has components given by :

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

with

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

being the temporal difference associated to transition $i \rightarrow j$, as introduced by Sutton [7]. When one uses a stochastic approximation of λ Policy Iteration, that is when the expectation $E_{\pi_{k+1}}$ is approximated by sampling, λ Policy Iteration reduces to the algorithm TD(λ) which is described in chapter 7 of [7]. Also, Bertsekas and Ioffe [1] showed that Approximate TD(λ) with a linear feature architecture, as described in chapter 8.2 of [7], corresponds to Approximate λ Policy Iteration where the value is updated by least square fitting using a gradient-type iteration after each sample. Eventually we do think that the ‘‘unified view’’ of Reinforcement Learning algorithms which is depicted in chapter 10.1 of [7] is in fact a picture of λ Policy Iteration.

Properties of Exact λ Policy Iteration In the original article introducing λ Policy Iteration [1], Bertsekas and Ioffe provide an analysis of the algorithm when it is run exactly (when $\epsilon_k = 0$). Define the following factor

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

We have $0 \leq \beta \leq \gamma < 1$. If $\lambda = 0$ (Value Iteration) then $\beta = \gamma$, and if $\lambda = 1$ (Policy Iteration) then $\beta = 0$.

They give some insight on what happens at each iteration and describe a constructive algorithm for estimating v_{k+1} from π_{k+1} and v_k :

Lemma 3 (One iterate of λ Policy Iteration[1])

For all k , define the operator

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

Assume $\mathcal{T}^{\pi_{k+1}}$ is a contraction mapping of modulus α for some norm $\|\cdot\|$. Then

- M_k is a contraction mapping of modulus $\beta\alpha$ for the same norm $\|\cdot\|$.
- The next iterate v_{k+1} is the (unique) fixed point of M_k .

The assumption of the lemma is always true with $\alpha = \gamma$ and the max-norm.

Then the authors show the convergence and provide an asymptotic rate of convergence :

⁵TD stands for Temporal Difference. The connection with TD algorithms was one of the motivation of the original article about λ Policy Iteration by Bertsekas and Ioffe [1]. Indeed, λ Policy Iteration is there also called ‘‘Temporal Difference Based Policy Iteration’’ and the presentation they give starts from the formulation of equation 8 (which is close to TD(λ) and then makes the connection with Value Iteration and Policy Iteration.

Theorem 4 (Convergence and Rate of Convergence of Exact λ Policy Iteration[1])

If the discount factor $\gamma < 1$, then v_k converges to v_* . Furthermore, after some index k_* , the rate of convergence is linear in β (for the norm of Lemma 3), that is

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

By making λ close to 1, β can be arbitrarily close to 0 so the above rate of convergence might look impressive. This needs to be put into perspective : the index k_* is the index after which the policy π_k does not change anymore (and is equal to the optimal policy π_*). As we said when we introduced of the algorithm, λ 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 λ , the less the operator M_k (Lemma 3) is contracting, and the more time it might take to compute the next iterate v_{k+1} (its fixed point).

Remark 1 Analyses are usually simpler for Value Iteration than for Policy Iteration. The key for studying Value Iteration is based on the fact that it computes the fixed point of the Bellman operator which is a γ -contraction mapping in max-norm. On the contrary, 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 ; Analyses of Policy Iteration rely on some other properties (like the fact that the sequence values is (approximately) non-decreasing). In fact, this observation can be drawn for λ Policy Iteration as soon as it does not reduce to Value Iteration :

As soon as $\lambda > 0$, there exists no norm for which the operator by which λ Policy Iteration updates the value from one iteration to the next is a contraction.

To see this, consider the following deterministic MDP with two states $\{1, 2\}$ and two actions $\{\text{change}, \text{stay}\}$: $r_1 = 0, r_2 = 1, P_{\text{change}}(s_2|s_1) = P_{\text{change}}(s_1|s_2) = P_{\text{stay}}(s_1|s_1) = P_{\text{stay}}(s_2|s_2) = 1$. Consider the following two value functions $v = (\epsilon, 0)$ and $v' = (0, \epsilon)$ with $\epsilon > 0$. Their corresponding greedy policies are $\pi = (\text{stay}, \text{change})$ and $\pi' = (\text{change}, \text{stay})$. Then, we can compute the next iterates of v and v' (using equation 9) :

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

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 all norms are equivalent, and as ϵ 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 $\lambda > 0$. \square

Approximate λ Policy Iteration There is a case study in the original article [1] describing an instance of Approximate λ Policy Iteration (with $\epsilon_k \neq 0$). However, to the best of our knowledge, there is no work studying the theoretical soundness of doing Approximate λ Policy Iteration.

3.5 Our contributions

Now that we have described the algorithms and some of their known properties, motivating the remaining of this paper is going to be straightforward. λ Policy Iteration is conceptually a very nice algorithm 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 all the results we have described so far to λ Policy Iteration (uniformly for all λ). The answer is positive :

- We provide a componentwise analysis of Exact and Approximate λ Policy Iteration.
- We show that the convergence rate of Exact λ Policy Iteration is (not only asymptotically) linear and we generalize the stopping criterion described for Value Iteration.
- We give (componentwise and L_p) bounds of the asymptotic error of Approximate λ Policy Iteration with respect to the *asymptotic* approximation error, Bellman residual, and Policy Bellman residual. We thus show that Approximate λ Policy Iteration is sound.
- We provide specific (better) bounds for the case when the value or the policy converges.
- Last but not least, we provide all our results using the span seminorms we have introduced in section 2, and using the relations between this span semi-norms and the standard L_p norms (equation 1), it can be seen that our results generalize all the previously described results.

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

4 Componentwise Performance bounds

This section contains our main results, which take the form of componentwise bounds. All the proofs are deferred to the appendices. The core of our work is the complete analysis of λ Policy Iteration (Appendix A). It serves as a basis for computing the rate of convergence of Exact λ Policy Iteration (Section 4.1, proof in Appendix B) and the asymptotic performance loss of Approximate λ Policy Iteration with respect to the approximation error (Section 4.2, proof in Appendix C). The asymptotic performance loss of Approximate λ Policy Iteration with respect to the Bellman residuals is somewhat simpler and is proved independently (Section 4.2, proof in Appendix D).

4.1 Performance bounds for Exact λ Policy Iteration

We first consider Exact λ Policy Iteration and provide some convergence rate bounds :

Lemma 4 (Componentwise Rate of Convergence bounds for Exact λ Policy Iteration)

The following matrices

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

are stochastic and

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

Remark 2 *These bounds imply that λ Policy Iteration converges to the optimal value function at least at a linear rate.*

4.2 Performance bounds for Approximate λ Policy Iteration

We provide componentwise bounds of the loss $v_* - v^{\pi_k}$ of using policy π_k instead of using the optimal policy, with respect to the approximation error ϵ_k , the Policy Bellman residual $\mathcal{T}_k v_k - v_k$ and the Bellman residual $\mathcal{T} v_k - v_k = \mathcal{T}_{k+1} v_k - v_k$. Note well the subtle difference between these two Bellman residuals : the Policy Bellman residual says how much v_k differs from the value of π_k while the Bellman residual says how much v_k differs from the value of the policies π_{k+1} and π_* .

Lemma 5 (Componentwise Performance bounds for Approximate λ Policy Iteration)

The following matrices

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

are stochastic and

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

Remark 3 *We can look at the relation between our bound for general λ and the bounds derived by Munos (Lemmas 1 and 2).*

- *Let us 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 we have just shown that $\limsup_{k \rightarrow \infty} v_ - v^{\pi_k}$ is upper bounded by :*

$$\limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} [(I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_{j+1} - (\gamma(I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_j + (P_*)^{k-j})] \epsilon_j. \quad (11)$$

The bound derived by Munos for Approximate Value Iteration (Lemma 1 page 4) is

$$\limsup_{k \rightarrow \infty} (I - \gamma P_k)^{-1} \sum_{j=0}^{k-1} \gamma^{k-j} [P_k P_{k-1} \dots P_{j+1} - (P_*)^{k-j}] \epsilon_j$$

$$\begin{aligned}
&= \limsup_{k \rightarrow \infty} \sum_{j=0}^{k-1} \gamma^{k-j} [(I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_{j+1} - (I - \gamma P_k)^{-1} (P_*)^{k-j}] \epsilon_j \\
&= \limsup_{k \rightarrow \infty} \sum_{j=0}^{k-1} \gamma^{k-j} [(I - \gamma P_k)^{-1} P_k P_{k-1} \dots P_{j+1} - ((I - \gamma P_k)^{-1} \gamma P_k (P_*)^{k-j} + (P_*)^{k-j})] \epsilon_j. \quad (12)
\end{aligned}$$

The above bounds are very close to each other : we go from equation 11 to equation 12 by replacing $P_{k-1} \dots P_j$ by $(P_*)^{k-j}$.

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

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

Then

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

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 2 page 6).

Roughly speaking, our componentwise analysis is thus a generalization of both analyses of Munos.

The bound 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.

Corollary 3 Suppose the value converges to some v . Write π its greedy policy and P the corresponding stochastic matrix. Define the following stochastic matrix :

$$B_v := (1 - \gamma) ((1 - \lambda)(I - \gamma P)^{-1} P + \lambda(I - \gamma P_*)^{-1} P)$$

Then the error necessarily converges to some ϵ and

$$v_* - v^\pi \leq \frac{\gamma}{1 - \gamma} [B_v - D] \epsilon$$

Corollary 4 *Suppose the policy converges to some π . Write P the corresponding stochastic matrix. The following matrices*

$$\begin{aligned} A^\pi &:= (1 - \lambda\gamma)P(I - \lambda\gamma P)^{-1} \\ A^\pi_{jk} &:= \frac{1 - \gamma}{\gamma^{k-j}} \left[\frac{\lambda\gamma}{1 - \lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} (P_*)^{k-1-i} (A^\pi)^{i-j} + \beta^{k-j} (I - \gamma P)^{-1} (A^\pi)^{k-1-j} \right] \\ B_{jk}^\pi &:= A^\pi_{jk} P \\ B'_{jk}^\pi &:= \frac{1 - \gamma}{1 - \lambda\gamma} \left[(P_*)^{k-j} + \frac{(1 - \lambda\gamma)\gamma(1 - \lambda)}{1 - \gamma} A^\pi_{jk} (I - \lambda\gamma P)^{-1} (P)^2 \right]. \end{aligned}$$

are stochastic and

$$v_* - v^\pi \leq \limsup_{k \rightarrow \infty} \frac{1 - \lambda\gamma}{1 - \gamma} \sum_{j=k_0}^{k-1} \gamma^{k-j} [B_{jk}^\pi - B'_{jk}^\pi] \epsilon_j$$

5 From componentwise bounds to norm and seminorm bounds

The componentwise bounds we have derived allow to derive L_p norm and seminorms bounds. To do so we will need two lemmas, the proofs of which are close to those of Munos [5, 4]; the main difference is that we come up with (better) span seminorm bounds instead of L_p /max-norms.

5.1 Deriving L_p norms and seminorms

Lemma 6 *Let x_k, y_k be sequences of vectors and X_k and X'_k sequences of stochastic matrices satisfying*

$$\limsup_{k \rightarrow \infty} |x_k| \leq K \limsup_{k \rightarrow \infty} (X_k - X'_k) y_k$$

For all distribution $\mu, \mu_k := \frac{1}{2}\mu(X_k + X'_k)$ is a distribution and

$$\limsup_{k \rightarrow \infty} \|x_k\|_{p, \mu} \leq K \limsup_{k \rightarrow \infty} \text{span}_{p, \mu_k} [y_k]$$

$$\limsup_{k \rightarrow \infty} \|x_k\|_\infty \leq K \limsup_{k \rightarrow \infty} \text{span}_\infty [y_k]$$

Proof: Write $a_{pk} := \arg \min_a \|y_k - ae\|_{p, \mu_k}$. As X_k and X'_k are stochastic matrices, $X_k e = X'_k e = e$, and we can write that :

$$\limsup_{k \rightarrow \infty} |x_k| \leq K \limsup_{k \rightarrow \infty} (X_k - X'_k)(y_k - a_{pk}e).$$

By taking the absolute value (we write $|x|$ the componentwise absolute value of x) we get

$$\limsup_{k \rightarrow \infty} |x_k| \leq K \limsup_{k \rightarrow \infty} (X_k + X'_k)|y_k - a_{pk}e|.$$

It can then be seen that

$$\begin{aligned} \limsup_{k \rightarrow \infty} \left(\|x_k\|_{p, \mu} \right)^p &= K^p \limsup_{k \rightarrow \infty} \mu(|x_k|)^p \\ &\leq K^p \limsup_{k \rightarrow \infty} \mu \left[\left(\frac{1}{2}(X_k + X'_k) \right) 2|y_k - a_{pk}e| \right]^p \end{aligned}$$

$$\begin{aligned}
&\leq K^p \limsup_{k \rightarrow \infty} \frac{1}{2} \mu(X_k + X'_k) (2|y_k - a_{pk}e|)^p \\
&= K^p \limsup_{k \rightarrow \infty} \mu_k (2|y_k - a_{pk}e|)^p \\
&= K^p \limsup_{k \rightarrow \infty} \left(2 \|y_k - a_{pk}e\|_{p, \mu_k} \right)^p \\
&= K^p \limsup_{k \rightarrow \infty} \left(\text{span}_{p, \mu_k} [y_k] \right)^p
\end{aligned}$$

where we used Jensen's inequality (using the convexity of $x \mapsto x^p$) and where the last inequality results from the definition of a_{pk} . \square

Lemma 7 Let x_k, y_k be sequences of vectors and X_{jk} and X'_{jk} sequences of stochastic matrices satisfying

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

For all distribution μ ,

$$\mu_{kj} := \frac{1}{2} \mu(X_{kj} + X'_{kj})$$

are distributions and

$$\limsup_{k \rightarrow \infty} \|x_k\|_{p, \mu} \leq \frac{K\gamma}{1-\gamma} \lim_{k_0 \rightarrow \infty} \sup_{k \geq j \geq k_0} \text{span}_{p, \mu_{kj}} [y_j].$$

$$\limsup_{k \rightarrow \infty} \|x_k\|_{\infty} \leq \frac{K\gamma}{1-\gamma} \limsup_{k \rightarrow \infty} \text{span}_{\infty} [y_k]$$

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

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

By taking the absolute value we get

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

It can then be seen that

$$\begin{aligned}
\limsup_{k \rightarrow \infty} \left(\|x_k\|_{p, \mu} \right)^p &= K^p \limsup_{k \rightarrow \infty} \mu (|x_k|)^p \\
&\leq K^p \limsup_{k \rightarrow \infty} \mu \left[\sum_{j=k_0}^{k-1} \gamma^{k-j} (X_{kj} + X'_{kj}) (|y_j - a_{pkj}e|) \right]^p \\
&= K^p \limsup_{k \rightarrow \infty} \mu \left[\frac{\left(\sum_{j=k_0}^{k-1} \gamma^{k-j} \frac{1}{2} (X_{kj} + X'_{kj}) 2 (|y_j - a_{pkj}e|) \right) \left(\sum_{j=k_0}^{k-1} \gamma^{k-j} \right)}{\sum_{j=k_0}^{k-1} \gamma^{k-j}} \right]^p \\
&\leq K^p \limsup_{k \rightarrow \infty} \mu \frac{\sum_{j=k_0}^{k-1} \gamma^{k-j} \frac{1}{2} (X_{kj} + X'_{kj}) \left[2 |y_j - a_{pkj}e| \left(\sum_{j=k_0}^{k-1} \gamma^{k-j} \right) \right]^p}{\sum_{j=k_0}^{k-1} \gamma^{k-j}}
\end{aligned}$$

$$\begin{aligned}
&= K^p \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} \mu_{kj} [2\|y_j - a_{pkj}e\|]^p \left(\sum_{j=k_0}^{k-1} \gamma^{k-j} \right)^{p-1} \\
&\leq K^p \left(\frac{\gamma}{1-\gamma} \right)^{p-1} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} [2\|y_j - a_{pkj}e\|_{p, \mu_{kj}}]^p \\
&= K^p \left(\frac{\gamma}{1-\gamma} \right)^{p-1} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} [\text{span}_{p, \mu_{kj}} [y_j]]^p \\
&\leq K^p \left(\frac{\gamma}{1-\gamma} \right)^{p-1} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} \left[\sup_{k' \geq j' \geq k_0} \text{span}_{p, \mu_{k'j'}} [y_{j'}] \right]^p \\
&= K^p \left(\frac{\gamma}{1-\gamma} \right)^{p-1} \frac{\gamma}{1-\gamma} \left[\sup_{k' \geq j' \geq k_0} \text{span}_{p, \mu_{k'j'}} [y_{j'}] \right]^p \\
&= K^p \left(\frac{\gamma}{1-\gamma} \right)^p \left[\sup_{k' \geq j' \geq k_0} \text{span}_{p, \mu_{k'j'}} [y_{j'}] \right]^p.
\end{aligned}$$

where x^p means the componentwise power of vector x , and where we used Jensen's inequality (with the convex function $x \mapsto x^p$) and the fact that $\sum_{j=k_0}^{k-1} \gamma^{k-j} \leq \frac{\gamma}{1-\gamma}$. 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. \square

5.2 Rates of convergence for Exact λ Policy Iteration

Theorem 5 (Rates of convergence for Exact λ Policy Iteration)

With the notations of Theorem 4, and for all distribution μ and indices $k > k_0$,

$$\begin{aligned}
\|v_* - v^{\pi_k}\|_{p, \mu} &\leq \frac{\gamma^{k-k_0}}{1-\gamma} \text{span}_{p, \frac{1}{2}\mu(F_{kk_0} + E'_{kk_0})} [v_* - v_{k_0}] \\
\|v_* - v^{\pi_k}\|_{\infty} &\leq \frac{\gamma^{k-k_0}}{1-\gamma} \text{span}_{\infty} [v_* - v_{k_0}] \\
\|v_* - v^{\pi_k}\|_{p, \mu} &\leq \frac{\gamma^{k-k_0}}{1-\gamma} \text{span}_{p, \frac{1}{2}\mu(E_{kk_0} + E'_{kk_0})} [\mathcal{T}v_{k_0} - v_{k_0}] \\
\|v_* - v^{\pi_k}\|_{\infty} &\leq \frac{\gamma^{k-k_0}}{1-\gamma} \text{span}_{\infty} [\mathcal{T}v_{k_0} - v_{k_0}] \\
\|v_* - v^{\pi_k}\|_{p, \mu} &\leq \gamma^{k-k_0} \left(\left\| (v_* - v_{k_0}) - \min_s [v_*(s) - v_{k_0}(s)]e \right\|_{p, \mu(P_*)^{k-k_0}} + \|v_*(s) - v^{\pi_{k_0+1}}\|_{\infty} \right) \\
\|v_* - v^{\pi_k}\|_{\infty} &\leq \gamma^{k-k_0} (\text{span}_{\infty} [v_* - v_{k_0}] + \|v_*(s) - v^{\pi_{k_0+1}}\|_{\infty})
\end{aligned}$$

The first pair of rate is expressed in terms of the distance between the value function and the optimal value function at some iteration k_0 . The second pair of inequalities can be used as a stopping criterion. Indeed, taking $k = k_0 + 1$ it implies for instance that

Theorem 6 (Stopping condition for Exact λ Policy Iteration)

If

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

then π_{k_0+1} is ϵ -optimal, that is $\|v_* - v^{\pi_{k_0+1}}\|_\infty \leq \epsilon$.

The last pair of inequalities rely 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 explicitly a $\frac{1}{1-\gamma}$ factor (though this factor is hidden in $v_* - v^{\pi_{k_0+1}}$).

5.3 Bounds for Approximate λ Policy Iteration

We can now derive the following bounds for λ Policy Iteration in span seminorm :

Theorem 7 *With the notations of theorem 5 and corollaries 3 and 4, and for all distribution μ*

$$\begin{aligned} \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{p,\mu} &\leq \frac{\gamma}{(1-\gamma)^2} \lim_{k_0 \rightarrow \infty} \sup_{k \geq j \geq k_0} \text{span}_{p, \frac{1}{2}\mu(B_{jk} + B'_{jk})} [\epsilon_j] \\ \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty &\leq \frac{\gamma}{(1-\gamma)^2} \limsup_{j \rightarrow \infty} \text{span}_\infty [\epsilon_j] \\ \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_{p,\mu} &\leq \frac{\gamma}{(1-\gamma)^2} \limsup_{k \rightarrow \infty} \text{span}_{p, \frac{1}{2}\mu(C_k + C'_k)} [\mathcal{T}_k v_k - v_k] \\ \limsup_{k \rightarrow \infty} \|v_* - v^{\pi_k}\|_\infty &\leq \frac{\gamma}{(1-\gamma)^2} \limsup_{k \rightarrow \infty} \text{span}_\infty [\mathcal{T}_k v_k - v_k] \\ \forall k, \|v_* - v^{\pi_k}\|_{p,\mu} &\leq \frac{\gamma}{1-\gamma} \text{span}_{p, \frac{1}{2}\mu(D + D'_k)} [\mathcal{T}v_{k-1} - v_{k-1}] \\ \forall k, \|v_* - v^{\pi_k}\|_\infty &\leq \frac{\gamma}{1-\gamma} \text{span}_\infty [\mathcal{T}v_{k-1} - v_{k-1}] \end{aligned}$$

If the value converges to some v , then the approximation error converges to some ϵ , and the corresponding greedy policy π satisfies :

$$\begin{aligned} \|v_* - v^\pi\|_{p,\mu} &\leq \frac{\gamma}{1-\gamma} \text{span}_{p, \frac{1}{2}\mu(B_v + D)} [\epsilon] \\ \|v_* - v^\pi\|_\infty &\leq \frac{\gamma}{1-\gamma} \text{span}_\infty [\epsilon] \end{aligned}$$

If the policy converges to some π , then :

$$\begin{aligned} \|v_* - v^\pi\|_{p,\mu} &\leq \frac{\gamma(1-\lambda\gamma)}{(1-\gamma)^2} \lim_{k_0 \rightarrow \infty} \sup_{k \geq j \geq k_0} \text{span}_{p, \frac{1}{2}\mu(B_{jk}^\pi + B'_{jk}^\pi)} [\epsilon_j] \\ \|v_* - v^\pi\|_\infty &\leq \frac{\gamma(1-\lambda\gamma)}{(1-\gamma)^2} \limsup_{j \rightarrow \infty} \text{span}_\infty [\epsilon_j] \end{aligned}$$

5.4 Using the concentration coefficients

We here consider the **concentration coefficient** introduced by Munos in [4, 5] and already mentioned before (equation 6 page 4). We recall its definition for clarity. We assume there exists a distribution ν and a real number $C(\nu)$ such that

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

Let X be an average of products of stochastic matrices of the MDP. Then,

$$\mu X y \leq C(\nu) \nu y$$

From this we can conclude analogues of lemmas 6 and 7.

Lemma 8 *Let x_k, y_k be sequences of vectors and X_k and X'_k sequences of stochastic matrices (that are averages of products of stochastic matrices of the MDP) satisfying*

$$\limsup_{k \rightarrow \infty} |x_k| \leq K \limsup_{k \rightarrow \infty} (X_k - X'_k) y_k$$

Then

$$\limsup_{k \rightarrow \infty} \|x_k\|_\infty \leq K [C(\nu)]^{1/p} \limsup_{k \rightarrow \infty} \text{span}_{p,\nu} [y_k]$$

Lemma 9 *Let x_k, y_k be sequences of vectors and X_{jk} and X'_{jk} sequences of stochastic matrices (that are averages of products of stochastic matrices of the MDP) satisfying*

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

Then

$$\limsup_{k \rightarrow \infty} \|x_k\|_\infty \leq \frac{K\gamma}{1-\gamma} [C(\nu)]^{1/p} \limsup_{k \rightarrow \infty} \text{span}_{p,\nu} [y_k]$$

As a consequence we can derive the following convergence rate bounds for Exact λ Policy Iteration.

Theorem 8

$$\|v_* - v^{\pi_k}\|_\infty \leq \frac{\gamma^{k-k_0}}{1-\gamma} [C(\nu)]^{1/p} \text{span}_{p,\nu} [v_* - v_{k_0}]$$

$$\|v_* - v^{\pi_k}\|_\infty \leq \frac{\gamma^{k-k_0}}{1-\gamma} [C(\nu)]^{1/p} \text{span}_{p,\nu} [\mathcal{T} v_{k_0} - v_{k_0}]$$

Also :

$$\|v_* - v^{\pi_k}\|_\infty \leq \gamma^{k-k_0} [C(\nu)]^{1/p} (\text{span}_{p,\mu} [v_* - v_{k_0}] + \|v_*(s) - v^{\pi_{k_0+1}}\|_\infty)$$

Also we have the following bounds for Approximate λ Policy Iteration :

Theorem 9

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

If the value converges to some v , then the approximation error converges to some ϵ , and the corresponding greedy policy π 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 π , then :

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

Références

- [1] D. Bertsekas and S. Ioffe. Temporal differences-based policy iteration and applications in neuro-dynamic programming. Technical Report LIDS-P-2349, MIT, 1996.
- [2] D.P. Bertsekas and J.N. Tsitsiklis. *Neurodynamic Programming*. Athena Scientific, 1996.
- [3] G. J. Gordon. Stable function approximation in dynamic programming. In Armand Prieditis and Stuart Russell, editors, *Proceedings of the Twelfth International Conference on Machine Learning*, pages 261–268, San Francisco, CA, 1995. Morgan Kaufmann.
- [4] R. Munos. Error bounds for approximate policy iteration. In *ICML*, pages 560–567, 2003.
- [5] R. Munos. Performance bounds in lp norm for approximate value iteration. *SIAM Journal on Control and Optimization*, 2007. To appear.
- [6] M. Puterman. *Markov Decision Processes*. Wiley, New York, 1994.
- [7] R.S. Sutton and A.G. Barto. *Reinforcement Learning, An introduction*. Bradford Book. The MIT Press, 1998.
- [8] R. Williams and L. Baird. Tight performance bounds on greedy policies based on imperfect value functions, 1993.

A General componentwise analysis of λ Policy Iteration

This section contains the core of all the remaining results. We show how to compute an upper bound of the loss for (approximate) λ Policy Iteration in general. It will be the basis for the derivation of componentwise bounds for approximate λ Policy Iteration (section 4.2) and exact λ Policy Iteration (section 4.1).

A.1 Overview of the analysis of the componentwise loss bound

We define :

- the **loss** of using policy π_k instead of the optimal policy :

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

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

$$w_k := v_k - \epsilon_k$$

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

$$d_k := v_* - w_k$$

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

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

- the **Bellman residual** between the $k - 1^{th}$ and the k^{th} values b.a. :

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

All our results come from a series of relations involving the above quantities. In this section, we will use the notation \bar{x} for an upper bound of x and \underline{x} for a lower bound.

Let us define the following stochastic matrix :

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

Then

Lemma 10 *The shift is related to the Bellman residual :*

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

Lemma 11 *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 :

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

Using Lemma 10, the bound on the Bellman residual also provides an upper on the shift :

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

Lemma 12 *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.

Then, an upper bound of the distance is :

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

Eventually, as

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

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

A.2 Proof of Lemma 10 : a relation between the shift $s_k = w_k - v^{\pi_k}$ and the Bellman residual $b_k = w_k - v_{k-1}$

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

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

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.3 Proof of Lemma 11 : an upper bound of the Bellman residual $b_k = T v_k - v_k$

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

$$\begin{aligned} b_k &= \mathcal{T}_{k+1} v_k - v_k \\ &= \mathcal{T}_{k+1} v_k - \mathcal{T}_k v_k + \mathcal{T}_k v_k - v_k \\ &\geq \mathcal{T}_k v_k - v_k \\ &= \mathcal{T}_k v_k - \mathcal{T}_k v^{\pi_k} + v^{\pi_k} - v_k \\ &= \gamma P_k (v_k - v^{\pi_k}) + v^{\pi_k} - v_k \\ &= (\gamma P_k - I)(s_k + \epsilon_k). \\ &= \beta A_k b_{k-1} + (\gamma P_k - I) \epsilon_k. \end{aligned} \tag{13}$$

where we eventually used the relation between s_k and b_k (Lemma 10). In other words :

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

with

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

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.4 Proof of Lemma 12 : an upper bound of the distance $d_k = v_* - w_k$

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 π^{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

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

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

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

B Componentwise rate of convergence for Exact λ Policy Iteration

We here derive the convergence rate bounds for Exact λ Policy Iteration (as expressed in Theorem 5 page 16). We rely on the loss bound analysis of section A with $\epsilon_k = 0$. In this specific case, we know that the loss $l_k \leq \overline{d}_k + \overline{s}_k$ where

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

We therefore have :

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

and

$$\overline{s}_k = \beta^{k-k_0} (I - \gamma P_k)^{-1} A_k A_{k-1} \dots A_{k_0+1} (-b_{k_0})$$

Therefore :

$$l_k \leq \left(\frac{\gamma^{k-k_0}}{1-\gamma} \right) E'_{kk_0} (-b_{k_0}) + \gamma^{k-k_0} (P_*)^{k-k_0} d_{k_0} \quad (14)$$

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 13 E'_{kk_0} is a stochastic matrix

Proof:

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

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

B.1 A bound with respect to the Bellman residual

We first need the following lemma :

Lemma 14 *The bias and the distance are related as follows :*

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

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

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

We thus have :

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

Then equation 14 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 14, we know that

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

Then equation 14 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

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

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 :

$$\begin{aligned} \hat{b}_{k_0} &\geq 0 \\ \mathcal{T}_{k_0+1} \hat{v}_{k_0} &\geq \hat{v}_{k_0} \end{aligned}$$

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

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 14 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 :

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

C Asymptotic componentwise loss bounds with respect to the approximation error

We here use the loss bound analysis of section A to derive an asymptotic analysis of approximate λ Policy Iteration with respect to the approximation error.

C.1 General analysis

Since

$$l_k = d_k + s_k \leq \overline{d}_k + \overline{s}_k, \tag{15}$$

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 :

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

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 :

$$\begin{aligned}
\overline{d}_k &= \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=k_0}^{k-1} \gamma^{k-1-i} \left(\sum_{j=k_0}^i \beta^{i-j} X_{i,j,k} (I - \gamma P_j) \epsilon_j + \mathcal{O}(\gamma^{i-k_0}) \right) - \sum_{i=k_0}^{k-1} \gamma^{k-i} (P_*)^{k-i} \epsilon_i + \mathcal{O}(\gamma^{k-k_0}) \\
&= \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=k_0}^{k-1} \sum_{j=k_0}^i \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} (I - \gamma P_j) \epsilon_j - \sum_{i=k_0}^{k-1} \gamma^{k-i} (P_*)^{k-i} \epsilon_i + \mathcal{O}(\gamma^{k-k_0}) \\
&= \frac{\lambda\gamma}{1-\lambda\gamma} \sum_{j=k_0}^{k-1} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} (I - \gamma P_j) \epsilon_j - \sum_{j=k_0}^{k-1} \gamma^{k-j} (P_*)^{k-j} \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \\
&= \sum_{j=k_0}^{k-1} \left[\left(\frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} (I - \gamma P_j) \right) - \gamma^{k-j} (P_*)^{k-j} \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \tag{16}
\end{aligned}$$

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

$$\begin{aligned}
\overline{s}_k &= \beta (I - \gamma P_k)^{-1} A_k (-\underline{b}_k) \\
&= \beta (I - \gamma P_k)^{-1} A_k \left[\left(\sum_{j=k_0}^{k-1} \beta^{k-1-j} (A_{k-1} A_{k-2} \dots A_{j+1}) (-x_j) \right) + \mathcal{O}(\gamma^{k-k_0}) \right] \\
&= \sum_{j=k_0}^{k-1} \frac{\beta^{k-j}}{1-\gamma} Y_{j,k} (I - \gamma P_j) \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \tag{17}
\end{aligned}$$

with

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

Eventually, from equations 15, 16 and 17 we get :

$$l_k \leq \sum_{j=k_0}^{k-1} \left[\left(\frac{\lambda\gamma}{1-\lambda\gamma} \sum_{i=j}^{k-1} \gamma^{k-1-i} \beta^{i-j} X_{i,j,k} + \frac{\beta^{k-j}}{1-\gamma} Y_{j,k} \right) (I - \gamma P_j) - \gamma^{k-j} (P_*)^{k-j} \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \tag{18}$$

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 15 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 norm is 1.

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

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

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

Equation 18 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 ,

$$\limsup_{k \rightarrow \infty} l_k \leq \frac{1}{1-\gamma} \limsup_{k \rightarrow \infty} \sum_{j=k_0}^{k-1} \gamma^{k-j} [B_{jk} - B'_{jk}] \epsilon_j \quad (19)$$

C.2 When the value converges

Suppose λ Policy Iteration converges to some value v . Let policy π 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 ϵ . Indeed from Algorithm 3 and equation 8, we get :

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

From the bound with respect to the Bellman residual (equation 24 page 30), 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

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

Lemma 16 B_v and D are stochastic matrices.

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

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

C.3 When the policy converges

Suppose λ Policy Iteration converges to some policy π . 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 15). 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

$$\begin{aligned}l_k &\leq \sum_{j=k_0}^{k-1} \left[\frac{\gamma^{k-j}}{1 - \gamma} A^\pi_{jk} [(1 - \lambda\gamma)P - \gamma(1 - \lambda)A^\pi P] - \gamma^{k-j} (P_*)^{k-j} \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \\ &= \sum_{j=k_0}^{k-1} \gamma^{k-j} \left[\frac{1 - \lambda\gamma}{1 - \gamma} A^\pi_{jk} P - \left[\frac{\gamma(1 - \lambda)}{1 - \gamma} A^\pi_{jk} A^\pi P + (P_*)^{k-j} \right] \right] \epsilon_j + \mathcal{O}(\gamma^{k-k_0}) \\ &= \frac{1 - \lambda\gamma}{1 - \gamma} \sum_{j=k_0}^{k-1} \gamma^{k-j} [B_{jk}^\pi - B'_{jk}{}^\pi] \epsilon_j + \mathcal{O}(\gamma^{k-k_0})\end{aligned}\tag{20}$$

where

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

Lemma 17 B_{jk}^π and $B'_{jk}{}^\pi$ are stochastic matrices.

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

$$\begin{aligned}\|B_{jk}^{\pi'}\| &= \frac{1-\gamma}{1-\lambda\gamma} \left(1 + \frac{\gamma(1-\lambda)}{1-\gamma}\right) \\ &= \frac{1-\gamma}{1-\lambda\gamma} \frac{1-\gamma+\gamma-\lambda\gamma}{1-\gamma} \\ &= 1. \quad \square\end{aligned}$$

D Componentwise bounds with respect to the Bellman residuals

In this section, 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 π_k while b_k says how much v_k differs from the value of the policies π_{k+1} and π_* .

D.1 Policy Bellman residual $b'_k = \mathcal{T}_k v_k - v_k$

Our analysis relies on the following lemma

Lemma 18 *Suppose we have a policy π , a function v that is an approximation of the value v^π of π in the sense that its residual $b' := \mathcal{T}^\pi v - v$ is small. Taking the greedy policy π' 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 π and π' .

Proof: We have :

$$\begin{aligned}v_* - v^{\pi'} &= \mathcal{T}_* v_* - \mathcal{T}^{\pi'} v^{\pi'} \\ &= \mathcal{T}_* v_* - \mathcal{T}_* v^\pi + \mathcal{T}_* v^\pi - \mathcal{T}_* v + \mathcal{T}_* v - \mathcal{T}^{\pi'} v + \mathcal{T}^{\pi'} v - \mathcal{T}^{\pi'} v^{\pi'} \\ &\leq \gamma P_*(v_* - v^\pi) + \gamma P_*(v^\pi - v) + \gamma P'(v - v^{\pi'})\end{aligned}\tag{21}$$

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

$$\begin{aligned}v^\pi - v &= \mathcal{T}^\pi v^\pi - v \\ &= \mathcal{T}^\pi v^\pi - \mathcal{T}^\pi v + \mathcal{T}^\pi v - v \\ &= \gamma P(v^\pi - v) + b' \\ &= (I - \gamma P)^{-1} b'\end{aligned}\tag{22}$$

and that

$$\begin{aligned}v - v^{\pi'} &= v - \mathcal{T}^{\pi'} v^{\pi'} \\ &= v - \mathcal{T}^\pi v + \mathcal{T}^\pi v - \mathcal{T}^{\pi'} v + \mathcal{T}^{\pi'} v - \mathcal{T}^{\pi'} v^{\pi'} \\ &\leq -b' + \gamma P'(v - v^{\pi'}) \\ &\leq (I - \gamma P')^{-1} (-b').\end{aligned}\tag{23}$$

where we used the fact that $\mathcal{T}^\pi v \leq \mathcal{T}^{\pi'} v$. We get the result by putting back equations 22 and 23 into equation 21. \square

To derive a bound for λ 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 Bellman residual $b_k = \mathcal{T} v_k - v_k$

We rely on the following lemma (which is for instance proved by Munos in [5])

Lemma 19 *Suppose we have a function v . Let π 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 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 22 we see that :

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

Thus

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

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

$$l_k \leq \frac{\gamma}{1 - \gamma} [D - D'_k] b_{k-1} \tag{24}$$

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.