

## A Relaxing assumptions

In this section we review the assumptions used throughout the paper and discuss if and how they could be relaxed.

### A.1 Assumption 1

We consider how to remove Asm. 1. When  $T_{\pi,t}(s) = 0$  we set  $\hat{\mu}_{\pi,t}(s)$  to an arbitrary default value  $\mu_\infty$ .<sup>7</sup> In this case, the prediction loss becomes

$$\mathcal{L}_n(\pi) = \frac{n}{S} \sum_{s \in \mathcal{S}} \mathbb{E}_\pi \left[ \frac{\sigma^2(s)}{T_{\pi,n}(s)} | T_{\pi,n}(s) > 0 \right] + E(\pi, n), \quad \text{with } E(\pi, n) := \frac{n}{S} \sum_{s \in \mathcal{S}} (\mu_\infty - \mu(s))^2 \mathbb{P}(T_{\pi,n}(s) = 0).$$

Asm. 1 makes the simplification that  $E(\pi, n) = 0$ . In order to deal with the general case, we need to take care of the event  $\{\exists s \in \mathcal{S}, T_{\pi,n}(s) = 0\}$  in which at least one state does not have any sample from which we could estimate its mean. An alternative is to consider that we initially have a “fictitious” observation equal to a fixed value at each state, which would introduce a small bias that tends to zero quickly. Another alternative could be to start by running a policy  $\pi_0$  over the states of the MDP and as soon as each state is visited at least once, we set the time step equal to 1 and begin our analysis. In the framework of the learning algorithm FW-AME, Asm. 1 can be easily replaced in practice by considering an adaptive length  $\tau_1$  such that at least one sample of each state is collected at the end of the first episode (which is what we do in the experiments in Sect. 5). The length of this phase would be small as the following result applies.

**Proposition 4.** *For any policy  $\pi \in \Pi^{SR}$ , under Asm. 2, the term  $E(\pi, n)$  decreases exponentially in  $n$ .*

*Proof.* Let  $n > 1/\eta_{\pi,\min}$ . Then setting  $\epsilon = \eta(s) - 1/n > 0$  yields<sup>8</sup>

$$\begin{aligned} \mathbb{P}(T_{\pi,n}(s) = 0) &= \mathbb{P}\left(T_{\pi,n}(s) < (\eta(s) - \epsilon)n\right) \\ &\leq \sqrt{\frac{2}{\eta_{\pi,\min}}} \exp\left(\frac{-n\gamma_{\text{ps}}^\pi \left(\eta(s) - \frac{1}{n}\right)^2}{16\eta(s)(1 - \eta(s))\left(1 + \frac{1}{n\gamma_{\text{ps}}^\pi}\right) + 40\left(\eta(s) - \frac{1}{n}\right)}\right). \end{aligned}$$

We thus obtain for any stationary policy  $\pi$  and any budget  $n > 1/\eta_{\pi,\min}$

$$E(\pi, n) \leq \frac{n}{S} \sum_{s \in \mathcal{S}} (\mu_\infty - \mu(s))^2 \sqrt{\frac{2}{\eta_{\pi,\min}}} \exp\left(\frac{-n\gamma_{\text{ps}}^\pi \left(\eta_{\pi,\min} - \frac{1}{n}\right)^2}{8\left(1 + \frac{1}{n\gamma_{\text{ps}}^\pi}\right) + 40\left(\eta_{\pi,\max} - \frac{1}{n}\right)}\right),$$

which proves the result. □

### A.2 Assumption 3

The reversibility assumption (Asm. 3) can be removed and Prop. 1, Lem. 1 as well as the proof of Thm. 1 could be easily adjusted to handle the case of non-reversible policies. As a result, the reversibility condition does not need to hold for the algorithm FW-AME and its vanishing regret guarantees. This can be achieved by replacing Prop. 1 with a concentration result adapted from Paulin et al. (2015).

**Proposition 5** (Thm. 3.10 and Prop. 3.14 from Paulin et al. (2015)). *Let us fix a stationary policy  $\pi$  which induces a time-homogeneous, ergodic Markov chain. We denote by  $P$  its transition matrix and by  $\hat{P}$  the time-reversal matrix of  $P$ . We denote by  $\eta_{\pi,\min} = \min_{s \in \mathcal{S}} \eta(s) > 0$  and  $\eta_{\pi,\max} = \max_{s \in \mathcal{S}} \eta(s)$  where  $\eta$  is the chain’s*

<sup>7</sup>Formally  $\mu_\infty = +\infty$  yet we can also set it equal to a suitable finite value depending on the distributions. For example, if the state distributions are Gaussian and the means belong to an interval  $[-\mu_{\max}, +\mu_{\max}]$ , we can fix  $\mu_\infty = 3\sigma_{\max} + \mu_{\max}$ , which ensures that the mean estimate computed from one single sample has an overwhelming probability of being more accurate than the default value  $\mu_\infty$  when there are no samples.

<sup>8</sup>Here we use the more general result for non-reversible chains reported in Prop. 5.

stationary distribution. We consider the pseudo-spectral gap  $\gamma_{ps}^\pi = \max_{k \geq 1} \gamma(\hat{P}^k P^k)/k > 0$ . For a given state  $s$  and for every  $\epsilon > 0$

$$\mathbb{P}\left(\left|\frac{T_{\pi,n}(s)}{n} - \eta(s)\right| > \epsilon\right) \leq \sqrt{\frac{2}{\eta_{\pi,\min}}} \exp\left(\frac{-n\gamma_{ps}^\pi \epsilon^2}{16\eta(s)(1-\eta(s))(1+1/(n\gamma_{ps}^\pi)) + 40\epsilon}\right).$$

In Sect. 4, the reversibility condition is intrinsically needed to relate the spectral gap with its spectral norm formulation, which is not possible for the pseudo-spectral gap. Nonetheless, rather than assuming that all policies are in the set of reversible stationary randomized policies  $\Pi^{\text{SRR}}$ , we could focus on computing a policy  $\pi_{\text{FMH}}^*$  belonging to the restricted set  $\Pi^{\text{SRR}}$ , thus replacing the assumption with an additional constrained in the optimization problem.

### A.3 Assumption 4

We can deal with the case when the MDP transition model  $p$  is unknown by following an optimistic approach similar to UCRL (Jaksch et al., 2010). We recall that the optimization problem solved by FW-AME at each episode is indeed equivalent to solving an MDP with known  $p$  and reward function set to  $\nabla \hat{\mathcal{L}}_{t_k-1}^+(\tilde{\lambda}_k)$ , which is already an optimistic evaluation of the true gradient. Whenever  $p$  is unknown, but an estimate and a confidence set are available, we can include the uncertainty of the estimate of  $p$  into the optimistic optimization of the MDP. Let us fix an episode  $k$  and  $t = t_k - 1$  the time step at the end of the previous episode. We introduce the following set that is  $p$ -dependent and thus unknown

$$\Lambda_{\underline{\eta}}^{(p)} = \left\{ \lambda \in \Delta(\mathcal{S} \times \mathcal{A}) : \forall s \in \mathcal{S}, \sum_{b \in \mathcal{A}} \lambda(s, b) \geq 2\underline{\eta} \quad \text{and} \quad \sum_{b \in \mathcal{A}} \lambda(s, b) = \sum_{s' \in \mathcal{S}, a \in \mathcal{A}} p(s|s', a) \lambda(s', a) \right\}.$$

The aim is to solve the following problem

$$\min_{\lambda \in \Lambda_{\underline{\eta}}^{(p)}} \min_{p \in \mathcal{C}_t} \langle c, \lambda \rangle = \sum_{s, a} \nabla \hat{\mathcal{L}}_t^+(\tilde{\lambda}_k)(s, a) \lambda(s, a).$$

If we define over  $\mathcal{S} \times \mathcal{A}$  the (bounded) reward function  $r = -\nabla \hat{\mathcal{L}}_t^+(\tilde{\lambda}_k)$ , we notice that the above problem can be reduced to the dual formulation of finding the policy that maximizes the average reward (Sect. 8, Puterman, 1994). As such, it becomes equivalent to solving the following problem

$$\max_{\pi \in \Lambda_{\underline{\eta}}^{(p)}} \max_{p \in \mathcal{C}_t} \rho_\pi(p), \tag{18}$$

where  $\rho_\pi(p)$  is the gain of stationary policy  $\pi$  in the MDP with transition probability function  $p$ . The confidence set  $\mathcal{C}_t$  defines a set of plausible transition probability functions at time  $t$ . Since the reward function is known, this corresponds to a set of plausible MDPs. Problem (18) thus returns an optimal policy in the plausible MDP with the largest gain. Lattimore and Szepesvári (Sect. 38, 2019) explicit the construction of  $\mathcal{C}_t$  and explain that the solutions of (18) are guaranteed to exist and can be found efficiently.<sup>9</sup>

While a complete derivation of the regret bound for this algorithm is left for future work, we expect the final result of Thm. 1 to remain unchanged. In fact, the optimal  $p$  returned by problem (18) belongs to  $\mathcal{C}_t$  so it is close to the real  $p$  up to a factor scaling in  $1/T_t$  by construction of  $\mathcal{C}_t$ . Hence, if the number of visits of any state-action pair (and not just the number of any state visit as in the case of known  $p$ ) is enforced to be proportional to the time step with high probability, then the derivation of the  $\tilde{O}(t^{-1/3})$  rate in the proof of Prop. 1 (cf. App. D) is unchanged.

<sup>9</sup>In a nutshell, the justification comes from introducing the extended Markov decision process  $\tilde{M}$  from Jaksch et al. (2010) and solving the average reward problem on that specific MDP using Extended Value Iteration. The fact that the extended action-sets of  $\tilde{M}$  are infinite is not problematic since  $\mathcal{C}_t$  is a convex polytope and has finitely many extremal points; as a result restricting the confidence sets to these points makes the extended MDP finite without changing the optimal policy.

---

**Algorithm 2** FMH
 

---

**Require:**  $\eta^*$  is the optimal stationary distribution of the convex problem (7).

**Require:** 3 parameters  $\rho_n$ ,  $\delta_n$  and  $\underline{\eta}$  (typically set respectively to  $S/n$ ,  $1/n$  and  $\min_s \eta^*(s)/2$ ).

Compute  $X_1$  the optimal solution of the convex problem ( $\mathcal{P}_1$ ) with parameters  $\rho_n$ ,  $\delta_n$  and  $\underline{\eta}$ .

Deduce the corresponding state distribution  $\eta_1$ :  $\eta_1(s) = \sum_{s' \in \mathcal{S}} X(s, s')$ .

Compute the optimal stationary policy of the convex problem ( $\mathcal{P}_2$ ) with  $\eta_1$  as target state distribution.

---

## B Faster-Mixing Heuristic FMH

### B.1 Derivation of the two-step method

FMH( $\pi^*$ ,  $n$ ) receives as input a budget  $n$  and  $\pi^*$ , the optimal solution of (7), and returns a stationary policy  $\pi_{\text{FMH}}^*$  by solving two convex optimization problems. An outline of FMH is provided in Alg. 2.

**Step 1.** In this step we first remove the stationarity constraint on  $\eta$  w.r.t. the MDP dynamics and replace it by a weaker but easier constraint involving the adjacency matrix of the MDP. Instead of using  $P$  as the kernel of the Markov chain associated to a policy, we consider it as a generic *transition matrix* that respects the possible transitions in the MDP, i.e.,  $P_{ij} = 0$  if  $Q_{ij} = 0$ . In this case, problem (11) becomes convex in  $P$  for a fixed  $\eta$  and convex in  $\eta$  for a fixed  $P$ , yet it is non-convex in both  $P$  and  $\eta$ . When  $P$  is fixed,  $\eta$  has no more degree of freedom (i.e., it can be directly derived from  $P$ ), thus any framework of alternate minimization cannot be applied here. We notice that the constraint of reversibility  $D_\eta P = P^T D_\eta$  is the toughest one to handle, since it involves both  $P$  and  $\eta$  and is not convex in  $P$  and  $\eta$ . This leads us to introduce the matrix variable  $X = D_\eta P \in \mathbb{R}^{S \times S}$  (i.e.,  $X_{ij} = \eta_i P_{ij}$ ). The reversibility constraint on  $P$  thus simply translates to a symmetry constraint on  $X$ . More discussion on the characteristics of the matrix  $X$  is for example provided in Hsu et al. (2015). We then obtain the following optimization problem with variable  $X$  (and its corresponding  $\eta$ )

$$\begin{aligned} \underset{X, \eta}{\text{minimize}} \quad & \mathcal{L}_0(X) := \sum_{s \in \mathcal{S}} \frac{\sigma^2(s)}{\eta(s)} + \rho_n \frac{1}{1 - \|D_\eta^{-1/2} X D_\eta^{-1/2} - \sqrt{\eta} \sqrt{\eta}^T\|_2} \\ \text{subject to} \quad & X \geq 0, \quad X = X^T, \quad \sum_{j \in \mathcal{S}} X_{ij} = \eta_i \quad \forall i \in \mathcal{S}, \quad X_{ij} = 0 \text{ if } Q_{ij} = 0, \quad \eta \geq \underline{\eta}, \quad \eta^T \mathbf{1} = 1. \end{aligned} \quad (19)$$

This problem is still non-convex in  $X$  and  $\eta$ . An idea could be to fix  $\eta$  and solve the convex problem in  $X$  (or equivalently  $P$ ). The most straightforward choice for  $\eta$  is to use  $\eta^*$ , the optimal stationary distribution of problem (7), and solve the convex problem of finding the fastest mixing Markov chain with stationary distribution  $\eta^*$  (Boyd et al., 2004). However the Markov chains whose stationary distributions are  $\eta^*$  might all mix poorly. Leveraging the intuition behind the regularized problem (11), we give more slack to  $\eta$  in order to find faster mixing Markov chains, at the cost of having  $\mathcal{L}(\eta)$  slightly larger than  $\mathcal{L}(\eta^*)$ , i.e., at the cost of a slightly worse asymptotic performance. We formalize this trade-off with the a parameter  $\delta_n$ , which represents how close we allow the stationary distribution  $\eta$  to be to  $\eta^*$  with respect to the  $\ell_2$ -norm (we pick the  $\ell_2$ -norm in order to ensure the convexity of the resulting constraint). We thus focus on solving the following surrogate optimization problem ( $\mathcal{P}_1$ )

$$\begin{aligned} \underset{X}{\text{minimize}} \quad & \mathcal{L}_1(X) := \sum_{s \in \mathcal{S}} \frac{\sigma^2(s)}{\sum_{j \in \mathcal{S}} X_{sj}} + \rho_n \frac{1}{1 - \|D_{\eta^*}^{-1/2} X D_{\eta^*}^{-1/2} - \sqrt{\eta^*} \sqrt{\eta^*}^T\|_2} \\ \text{subject to} \quad & X \geq 0, \quad X = X^T, \quad X_{ij} = 0 \text{ if } Q_{ij} = 0, \\ & \|D_{\eta^*}^{-1/2} X D_{\eta^*}^{-1/2} - \sqrt{\eta^*} \sqrt{\eta^*}^T\|_2 \leq 1, \\ & \sum_{(i,j) \in \mathcal{S}^2} X_{ij} = 1, \quad \sum_{j \in \mathcal{S}} X_{ij} \geq \underline{\eta}, \quad \sum_{i \in \mathcal{S}} \left( \sum_{j \in \mathcal{S}} X_{ij} - \eta_i^* \right)^2 \leq \delta_n^2, \end{aligned} \quad (\mathcal{P}_1)$$

where the small positive constant  $\underline{\eta}$  should satisfy  $\underline{\eta} \leq \min_s \eta^*(s)$ . Prop. 6 guarantees the convexity and feasibility of the optimization problem ( $\mathcal{P}_1$ ).

**Proposition 6.** ( $\mathcal{P}_1$ ) is convex in  $X$  and well-defined for any  $\delta_n$ .

*Proof.* The convexity of  $(\mathcal{P}_1)$  is easily obtained from the convexity of the non-regularized problem, the convexity of the function  $X \mapsto 1/(1 - \|X\|_2)$  and the convexity of all the constraints. There can however exist some matrices  $X$  such that  $\|D_{\eta^*}^{-1/2} X D_{\eta^*}^{-1/2} - \sqrt{\eta^*} \sqrt{\eta^*}^T\|_2 \geq 1$ , thus making  $(\mathcal{P}_1)$  either undefined in its objective function (if the norm is equal to 1) or not satisfying one of the constraints. We thus need to ensure that for any fixed  $\delta_n$  there exists at least one matrix  $X$  such that  $\|D_{\eta^*}^{-1/2} X D_{\eta^*}^{-1/2} - \sqrt{\eta^*} \sqrt{\eta^*}^T\|_2 < 1$  with the remaining constraints satisfied. To do so, we introduce the transition matrix  $M^* = (P^* + \widehat{P}^*)/2$  with  $\widehat{P}^*$  the time-reversed transition matrix of  $P^*$  which is the transition matrix of the optimal policy for problem (7). Whereas  $P^*$  is not necessarily reversible w.r.t.  $\eta^*$ , it is the case for  $M^*$ , thus yielding  $\text{SLEM}(M^*) < 1$ . We also define  $X = D_{\eta^*} M^*$ . By construction of  $X$ , we have  $X \geq 0$ ,  $X = X^T$ ,  $\sum_{i,j} X_{ij} = 1$ ,  $X_{ij} = 0$  if  $Q_{ij} = 0$  and  $\sum_j X_{ij} \geq \underline{\eta}$ . Furthermore, we have  $\sum_i (\sum_j X_{ij} - \eta_i^*)^2 = \sum_i (\sum_j \eta_i^* M_{ij}^* - \eta_i^*)^2 = \sum_i (\eta_i^*)^2 (\sum_j M_{ij}^* - 1)^2 = 0$  which means that all the constraints are verified. In addition, since  $M^*$  is reversible w.r.t.  $\eta^*$ , we have  $\|D_{\eta^*}^{-1/2} X D_{\eta^*}^{-1/2} - \sqrt{\eta^*} \sqrt{\eta^*}^T\|_2 = \|D_{\eta^*}^{1/2} M^* D_{\eta^*}^{-1/2} - \sqrt{\eta^*} \sqrt{\eta^*}^T\|_2 = \text{SLEM}(M^*) < 1$ . This proves that  $(\mathcal{P}_1)$  is well-defined.  $\square$

Solving the convex optimization problem  $(\mathcal{P}_1)$  yields an optimal matrix  $X_1 \in \mathbb{R}^{S \times S}$ , from which we easily obtain the associated stationary distribution  $\eta_1$  as well as the transition matrix of the associated Markov chain  $P_1$

$$\eta_1(s) = \sum_{s' \in \mathcal{S}} X_1(s, s') \quad \text{and} \quad P_1(s, s') = \frac{X_1(s, s')}{\sum_{s'} X_1(s, s')}.$$

**Step 2.** The distribution  $\eta_1$  is stationary w.r.t. the Markov chain  $P_1$  (which is expected to have better mixing properties than  $P^*$ ), but it may not be feasible w.r.t. the MDP dynamics. As a result, we must now find a stationary policy  $\pi$  whose stationary distribution is closest to  $\eta_1$ . This is closely linked to the steady-state control problem from Akshay et al. (2013), where it is proved that for an ergodic MDP the problem of finding a stationary policy given a target stationary state distribution is effectively decidable in polynomial time. If the steady-state control problem admits a solution, such a policy can be computed by simply solving a polynomial-size linear program. More precisely, we seek a policy  $\pi$  in the set of non-negative reals  $\{\pi_{s,a} | s \in \mathcal{S}, a \in \mathcal{A}_s\}$  such that

$$\forall s \in \mathcal{S}, \quad \sum_{s' \in \mathcal{S}, a \in \mathcal{A}_{s'}} \eta_1(s') p(s|s', a) \pi_{s',a} = \eta_1(s) \quad \text{and} \quad \sum_{a \in \mathcal{A}_s} \pi_{s,a} = 1.$$

If the steady-state control problem does not admit a solution, we seek a stationary policy whose stationary state distribution is closest to  $\eta_1$  w.r.t. the  $\ell^2$ -norm by solving the following convex optimization problem  $(\mathcal{P}_2)$  in  $\pi$

$$\begin{aligned} & \underset{\pi}{\text{minimize}} \quad \sum_s \left( \eta_1(s) - \sum_{s' \in \mathcal{S}, a \in \mathcal{A}_{s'}} \eta_1(s') p(s|s', a) \pi_{s',a} \right)^2 & (\mathcal{P}_2) \\ & \text{subject to} \quad \forall s \in \mathcal{S}, \quad \pi_{s,a} \geq 0 \quad \forall a \in \mathcal{A}_s \quad \text{and} \quad \sum_{a \in \mathcal{A}_s} \pi_{s,a} = 1. \end{aligned}$$

Since we do not know in advance if the steady-state control problem admits a solution or not, we directly solve problem  $(\mathcal{P}_2)$  which encompasses both cases (its optimal value is 0 if the steady-state control problem admits a solution). This yields a policy denoted  $\pi_{\text{FMH}}^*$ .

## B.2 A SDP formulation of FMH (FMH-SDP)

We notice that step 1 of FMH is by far the most computationally demanding, due to the complexity of the objective function and constrained set of problem  $(\mathcal{P}_1)$ . Fortunately, the symmetry constraint on  $X$  leads to the symmetry of the matrix  $D_{\eta^*}^{-1/2} X D_{\eta^*}^{-1/2} - \sqrt{\eta^*} \sqrt{\eta^*}^T$ , which is a very useful property because it becomes easy to compute a subgradient of its spectral norm w.r.t.  $X$  (see e.g., Boyd et al. (Sect. 5.1, 2004)). We can thus apply subgradient descent to solve  $(\mathcal{P}_1)$ . However a projection on the constrained set is required at each step. We thus propose an alternative method to solve problem  $(\mathcal{P}_1)$  that is projection-free and hence more computationally efficient. Since this approach uses semidefinite programming, the resulting heuristic is called FMH-SDP.

The key observation is that the regularizing term in  $(\mathcal{P}_1)$  partially “takes into account” the non-regularized one through the last constraint  $\|\eta - \eta^*\| \leq \delta_n$ . Furthermore, the regularizing term corresponds (up to composition

of a non-decreasing function) to minimizing the spectral norm of a symmetric matrix. Drawing inspiration from Boyd et al. (Sect. 2.3, 2004), we can express it as a semidefinite program (SDP) which can be solved efficiently using standard SDP solvers. Introducing a scalar variable  $s$  to bound the spectral norm, step 1 of FMH is replaced by the following SDP problem whose variables are the matrix  $X$  and the scalar  $s$

$$\begin{aligned}
 & \underset{X, s}{\text{minimize}} && s \\
 & \text{subject to} && -sI \preceq D_{\eta^*}^{-1/2} X D_{\eta^*}^{-1/2} - \sqrt{\eta^*} \sqrt{\eta^*}^T \preceq sI \\
 & && X \geq 0, \quad X = X^T, \quad X_{ij} = 0 \text{ if } J_{ij} = 0 \\
 & && \sum_{(i,j) \in \mathcal{S}^2} X_{ij} = 1, \quad \sum_{j \in \mathcal{S}} X_{ij} \geq \underline{\eta}, \quad |(\sum_{j \in \mathcal{S}} X_{ij})_i - \eta_i^*| \leq (\delta_n)_i.
 \end{aligned} \tag{20}$$

FMH-SDP is not only more computationally efficient due to its SDP formulation but it also loses the dependency on the hyper-parameter  $\rho_n$  as only  $\delta_n$  remains.

## C Proofs

We first recall the performance loss suffered by the continuous relaxation in the bandit case, where the frequency  $T_{\pi, n}/n$  is replaced by an allocation  $\lambda$  in the simplex. In order to keep the notation as consistent as possible, consider a stochastic bandit problem with  $S$  arms, let  $\Delta_n = \{\eta_n \in [0, 1]^S : \eta_n(s) = \frac{T_n(s)}{n}\}$  and  $\Delta$  be the discrete and continuous simplex over  $S$  arms, where  $\eta_n(s)$  is the frequency associated to  $T_n$  pulls. Since in this case a policy directly selects arms rather than actions, the objective functions  $\mathcal{L}_n$  and  $\mathcal{L}$  coincide and we can write

$$\mathcal{L}(\eta) = \frac{1}{S} \sum_s \frac{\sigma^2(s)}{\eta(s)},$$

where  $\eta$  may be either a discrete or a continuous allocation. We have the following.

**Proposition 7.** *Let  $\eta_n^* = \arg \min_{\eta_n \in \Delta_n} \mathcal{L}(\eta_n)$  be the optimal discrete allocation. As computing  $\eta_n^*$  is NP hard, a standard solution is to first compute  $\eta^* = \arg \min_{\eta \in \Delta} \mathcal{L}(\eta)$  and then round it to obtain  $\tilde{\eta}_n$ . If  $\tilde{\eta}_n$  is computed using efficient apportionment techniques (Chapter 12, Pukelsheim, 2006), then for any budget  $n > 2S$  we have*

$$\mathcal{L}(\tilde{\eta}_n) - \mathcal{L}(\eta_n^*) \leq \frac{2}{n} \sum_s \frac{\sigma^2(s)}{\eta^*(s)} = \frac{2S}{n} \mathcal{L}(\eta^*).$$

Furthermore, for any  $n \geq 4/(S\eta_{\min}^2)$ , where  $\eta_{\min} = \min_s \eta^*(s)$  we have

$$\mathcal{L}(\tilde{\eta}_n) - \mathcal{L}(\eta_n^*) \leq \frac{8\sigma_{\max}^2}{\eta_{\min}^3 n^2}.$$

*Proof.* Using efficient apportionment techniques for rounding we have (Lem. 12.8, Pukelsheim, 2006)

$$\min_s \frac{\tilde{\eta}_n(s)}{\eta^*(s)} \geq 1 - \frac{S}{n}, \quad \text{i.e.,} \quad \forall s, \tilde{\eta}_n(s) \geq \eta^*(s) \left(1 - \frac{S}{n}\right),$$

which also implies the other direction as

$$\tilde{\eta}_n(s) = 1 - \sum_{s' \neq s} \tilde{\eta}_n(s') \leq 1 - \sum_{s' \neq s} \eta^*(s') + \sum_{s' \neq s} \eta^*(s') \frac{S}{n} \leq \eta^*(s) + \frac{S}{n}.$$

Then we can bound the performance loss of  $\tilde{\eta}_n$  as

$$\mathcal{L}(\tilde{\eta}_n) - \mathcal{L}(\eta_n^*) = \mathcal{L}(\tilde{\eta}_n) - \mathcal{L}(\eta^*) + \underbrace{\mathcal{L}(\eta^*) - \mathcal{L}(\eta_n^*)}_{\leq 0} \leq \frac{1}{S} \sum_s \sigma^2(s) \left( \frac{1}{\tilde{\eta}_n(s)} - \frac{1}{\eta^*(s)} \right).$$

Under the assumption that  $n > 2S$ , we can bound each of the summands as

$$\frac{1}{\tilde{\eta}_n(s)} - \frac{1}{\eta^*(s)} = \frac{\eta^*(s) - \tilde{\eta}_n(s)}{\tilde{\eta}_n(s)\eta^*(s)} \leq \frac{\eta^*(s)S/n}{(\eta^*(s))^2(1 - S/n)} \leq \frac{2S}{\eta^*(s)n},$$

which proves the  $O(1/n)$  upper bound. Recalling the definition of  $\mathcal{L}(\eta^*)$  we obtain the final statement

$$\mathcal{L}(\tilde{\eta}_n) - \mathcal{L}(\eta_n^*) \leq 2 \sum_s \frac{\sigma^2(s)}{\eta^*(s)n} = \frac{2S}{n} \mathcal{L}(\eta^*).$$

An even faster rate can be obtained exploiting the smoothness of  $\mathcal{L}$ . Let  $\bar{\Delta} = \{\eta \in \Delta : \forall s, \eta(s) \geq \eta_{\min}/2\}$ . Since  $\tilde{\eta}(s) \geq \eta^*(s)(1 - S/n)$ , for any any  $n > 2S$  we have  $\tilde{\eta}_n, \eta^* \in \bar{\Delta}$ , hence using the  $\bar{C}$ -smoothness of  $\mathcal{L}$  on  $\bar{\Delta}$  with  $\bar{C} = \frac{2\sigma_{\max}^2}{S(\eta_{\min}/2)^3}$ , we can write (Thm. 12.10 Pukelsheim, 2006)

$$\mathcal{L}(\tilde{\eta}_n) - \mathcal{L}(\eta_n^*) \leq \frac{\bar{C}}{2} \|\tilde{\eta}_n - \eta^*\|_2^2 \leq \frac{8\sigma_{\max}^2}{\eta_{\min}^3 n^2},$$

which corresponds to an asymptotic rate of  $O(1/n^2)$ . Since the multiplicative constants are larger than those of the  $O(1/n)$  rate, the rate  $O(1/n^2)$  effectively starts when  $n$  is big enough. A rough bound on  $n$  for the second bound to be effectively smaller than the first is obtained by upper-bounding  $\mathcal{L}(\eta^*) \leq \sigma_{\max}^2/\eta_{\min}$  as

$$\frac{8\sigma_{\max}^2}{\eta_{\min}^3 n^2} \leq \frac{2\sigma_{\max}^2 S}{\eta_{\min} n} \iff n \geq \frac{4}{S\eta_{\min}^2},$$

which concludes the proof. □

*Proof of Proposition 1.* The first statement is a direct application of the relationship between mixing and spectral gap. For any policy reversible and ergodic policy  $\pi$ , any starting state  $s'$  and any state  $s$ , we have from Diaconis et al. (Prop. 3, 1991)

$$|\mathbb{P}_\pi(s_t = s | s_1 = s') - \eta_\pi(s)| \leq \frac{1}{2} \sqrt{\frac{1 - \eta_\pi(s')}{\eta_\pi(s')}} (1 - \gamma_\pi)^t$$

Then the difference between the expected frequency and the stationary distribution is bounded as

$$\begin{aligned} \left| \frac{\mathbb{E}[T_{\pi,n}(s)]}{n} - \eta_\pi(s) \right| &\leq \frac{1}{n} \sum_{t=2}^n |\mathbb{P}_\pi(s_t = s | s_1 = \bar{s}) - \eta_\pi(s)| \\ &\leq \frac{1}{2\sqrt{\eta_{\min} n}} \sum_{t=1}^n (1 - \gamma_\pi)^t \leq \frac{1}{2\sqrt{\eta_{\min} n} \gamma_\pi}. \end{aligned}$$

□

*Proof of Proposition 3.* The state-action polytope  $\Lambda$  is closed, bounded and convex according to Puterman (Thm. 8.9.4, 1994). The problem (7) is thus convex in  $\lambda$  due to the convexity of the objective function and constraints. It is straightforward that  $\pi_{\lambda^*} \in \Pi^{\text{SR}}$ . From Puterman (Thm. 8.8.1, 1994), the stationary distribution  $\eta_{\pi_{\lambda^*}}$  of  $\pi_{\lambda^*}$  is the unique solution of the system of equations  $\sum_{s'} P_{\pi_{\lambda^*}}(s|s') \eta_{\pi_{\lambda^*}}(s') = \eta_{\pi_{\lambda^*}}(s)$  (for each state  $s$ ) subject to  $\sum_s \eta_{\pi_{\lambda^*}}(s) = 1$ . Given that  $(\sum_a \lambda^*(s, a))_s$  is a solution, it corresponds to the stationary distribution  $\eta_{\pi_{\lambda^*}}$ . By contradiction, assume that there exists a policy  $\bar{\pi} \in \Pi^{\text{SR}}$  such that  $\mathcal{L}(\bar{\pi}, \eta_{\bar{\pi}}) < \mathcal{L}(\pi_{\lambda^*}, \eta_{\pi_{\lambda^*}})$ . Then define for every state-action pair  $(s, a)$  the quantity  $\bar{\lambda}(s, a) = \eta_{\bar{\pi}}(s) \bar{\pi}(a|s)$ . It is evident that  $\bar{\lambda} \in \Delta(\mathcal{S} \times \mathcal{A})$ , furthermore for every state  $s$ , we have

$$\begin{aligned} \sum_{s', a} p(s|s', a) \bar{\lambda}(s', a) &= \sum_{s', a} p(s|s', a) \eta_{\bar{\pi}}(s') \bar{\pi}(a|s) = \sum_{s'} \eta_{\bar{\pi}}(s') \sum_a p(s|s', a) \bar{\pi}(a|s') \\ &= \sum_{s'} \eta_{\bar{\pi}}(s') P_{\bar{\pi}}(s|s') = \eta_{\bar{\pi}}(s) = \sum_a \bar{\lambda}(s, a), \end{aligned}$$

since by stationarity of the policy  $\bar{\pi}$ , the Markov chain transition matrix  $P_{\bar{\pi}}$  is stationary w.r.t.  $\eta_{\bar{\pi}}$ . So  $\bar{\lambda}$  satisfies the constraint of stationarity of (7), and

$$\mathcal{L}(\bar{\lambda}) = \sum_s \frac{\sigma^2(s)}{\eta_{\bar{\pi}}(s)} = \mathcal{L}(\bar{\pi}, \eta_{\bar{\pi}}) < \mathcal{L}(\pi_{\lambda^*}, \eta_{\pi_{\lambda^*}}) = \sum_s \frac{\sigma^2(s)}{\eta_{\pi_{\lambda^*}}(s)} = \sum_s \frac{\sigma^2(s)}{\sum_a \lambda^*(s, a)} = \mathcal{L}(\lambda^*),$$

which contradicts the optimality of  $\lambda^*$  for problem (7) and thus proves that  $\pi_{\lambda^*}$  is the optimal solution of the problem (8). Finally, the upper bound on the smoothness parameter  $C_{\underline{\eta}}$  on the restricted set  $\Lambda_{\underline{\eta}}$  is derived using that the maximal eigenvalue of a symmetric block matrix with positive eigenvalues is bounded by above by the sum of maximal eigenvalues of its diagonal blocks.  $\square$

*Proof of Lemma 1.* The proof is a rather direct application of Prop. 1. We first recall the exact formulation of the term  $\epsilon_{\pi}(s, n, \delta)$  in Prop. 1 (see e.g., Hsu et al. (Thm. 3, 2015), Paulin et al. (Thm. 3.8, 2015)):

$$\epsilon_{\pi}(s, n, \delta) := \sqrt{8\eta_{\pi}(s)(1 - \eta_{\pi}(s)) \frac{\ln(\frac{1}{\delta} \sqrt{\frac{2}{\eta_{\pi, \min}}})}{\gamma_{\pi} n} + 20 \frac{\ln(\frac{1}{\delta} \sqrt{\frac{2}{\eta_{\pi, \min}}})}{\gamma_{\pi} n}}.$$

Let  $\eta_{\pi, n}(s) = \frac{T_{\pi, n}(s)}{n}$  be the empirical frequency of visits to state  $s$ . Since we need all following statements to hold simultaneously for all states  $s \in \mathcal{S}$  and all stationary policies  $\pi \in \Pi^{\text{SR}}$ , we need to take a union bound over states and a cover over the action simplex at each state, which leads to tuning  $\delta = \delta'/(SA^S)$  in the high-probability guarantees of Prop. 1, which then hold with probability  $1 - \delta'$ . Furthermore, we have the following deterministic bound

$$\left| \frac{1}{\eta_{\pi, n}(s)} - \frac{1}{\eta_{\pi}(s)} \right| \leq \max\{n, \frac{1}{\eta_{\pi}(s)}\},$$

where we used Asm. 1 to ensure that  $1/\eta_{\pi, n}(s) \leq n$ . We introduce the event

$$\mathcal{E}_1(s, n, \delta) = \{\eta_{\pi, n}(s) \geq \eta_{\pi}(s) - \epsilon_{\pi}(s, n, \delta)\}.$$

Then we have

$$\begin{aligned} \left| \mathbb{E} \left[ \frac{1}{\eta_{\pi, n}(s)} - \frac{1}{\eta_{\pi}(s)} \right] \right| &\leq \left| \mathbb{E} \left[ \left( \frac{1}{\eta_{\pi, n}(s)} - \frac{1}{\eta_{\pi}(s)} \right) \mathbb{I}\{\mathcal{E}_1(s, n, \delta)\} \right] \right| + \left| \mathbb{E} \left[ \left( \frac{1}{\eta_{\pi, n}(s)} - \frac{1}{\eta_{\pi}(s)} \right) \mathbb{I}\{\mathcal{E}_1^c(s, n, \delta)\} \right] \right| \\ &\leq \left| \mathbb{E} \left[ \left( \frac{1}{\eta_{\pi, n}(s)} - \frac{1}{\eta_{\pi}(s)} \right) \mathbb{I}\{\mathcal{E}_1(s, n, \delta)\} \right] \right| + \max\{n, \frac{1}{\eta_{\pi}(s)}\} \mathbb{P}\{\mathcal{E}_1^c(s, n, \delta)\} \\ &\leq \left| \mathbb{E} \left[ \frac{\eta_{\pi}(s) - \eta_{\pi, n}(s)}{\eta_{\pi, n}(s)\eta_{\pi}(s)} \mathbb{I}\{\mathcal{E}_1(s, n, \delta)\} \right] \right| + \max\{n, \frac{1}{\eta_{\pi}(s)}\} \delta' \\ &\leq \frac{|\mathbb{E}[\eta_{\pi}(s) - \eta_{\pi, n}(s)]|}{\eta_{\pi}(s)(\eta_{\pi}(s) - \epsilon_{\pi}(s, n, \delta))} + \max\{n, \frac{1}{\eta_{\pi}(s)}\} \delta' \\ &\leq \frac{1}{2\sqrt{\eta_{\min}} n \gamma_{\pi} \eta_{\pi}^2(s)} \left( 1 + 2 \frac{\epsilon_{\pi}(s, n, \delta)}{\eta_{\pi}(s)} \right) + \max\{n, \frac{1}{\eta_{\pi}(s)}\} \delta', \end{aligned}$$

where the last inequality follows from  $1/(1-x) \leq 1+2x$  for  $0 < x \leq 1/2$  which can be applied due to the condition that  $n$  is big enough so that  $\epsilon_{\pi}(s, n, \delta) \leq \eta_{\pi}(s)/2$ . Since this condition requires  $n \geq O(1/\eta_{\min}^2)$ , we can resolve the maximum in the previous expression as  $\max\{n, \frac{1}{\eta_{\pi}(s)}\} \leq n$ . Finally, setting  $\delta' = 1/n^2$  translates to the inequality on the objective function

$$|\mathcal{L}_n(\pi) - \mathcal{L}(\pi, \eta_{\pi})| \leq \ell_n(\pi) := \frac{1}{S\sqrt{\eta_{\min}} n \gamma_{\pi}} \sum_{s \in \mathcal{S}} \frac{\sigma^2(s)}{\eta_{\pi}^2(s)} \left( 1 + 2 \frac{\epsilon_{\pi}(s, n, \delta)}{\eta_{\pi}(s)} \right),$$

from which we obtain the final statement as

$$\begin{aligned} \mathcal{L}_n(\pi_{\lambda^*}) - \mathcal{L}_n(\pi_n^*) &\leq \mathcal{L}(\pi_{\lambda^*}, \eta_{\pi_{\lambda^*}}) + \ell_n(\pi_{\lambda^*}) - \mathcal{L}(\pi_n^*, \eta_{\pi_n^*}) + \ell_n(\pi_n^*) \\ &\leq \ell_n(\pi_{\lambda^*}) + \ell_n(\pi_n^*). \end{aligned}$$

$\square$

*Proof of Lemma 2.* The proof relies on the concentration inequality in Eq. 9. We proceed through the following inequalities

$$\mathcal{L}_n(\pi_{\text{reg}}^*) \leq \mathcal{L}(\pi_{\text{reg}}^*, \eta_{\pi_{\text{reg}}^*}) + \ell_n(\pi_{\text{reg}}^*) \leq \mathcal{L}(\pi_n^*, \eta_{\pi_n^*}) + \ell_n(\pi_n^*) \leq \mathcal{L}_n(\pi_n^*) + 2\ell_n(\pi_n^*),$$

where in the first and last inequality we used Eq. 9, and where the second inequality follows from the definition of  $\pi_{\text{reg}}^*$  as the optimal solution to the regularized problem.  $\square$

*Proof of Lemma 3.* Introducing the term  $H := \mathcal{L}^{\text{reg}}(\pi_{\text{FMH}}^*) - \mathcal{L}^{\text{reg}}(\pi_{\text{reg}}^*)$  where  $\mathcal{L}^{\text{reg}}$  is defined in Eq. 11, we have

$$\mathcal{L}_n(\pi_{\text{FMH}}^*) \leq \mathcal{L}^{\text{reg}}(\pi_{\text{FMH}}^*) = H + \mathcal{L}^{\text{reg}}(\pi_{\text{reg}}^*) \leq H + \mathcal{L}_n(\pi_n^*) + 2\ell_n(\pi_n^*).$$

Given the expression of  $\ell_n(\pi)$  provided in Lem. 1, we can write

$$\mathcal{L}^{\text{reg}}(\pi, \eta) = \mathcal{L}(\pi, \eta) + \frac{\rho_n}{1 - \|D_\eta^{1/2} P_\pi D_\eta^{-1/2} - \sqrt{\eta} \sqrt{\eta}^\top\|_2} + O(n^{-3/2}).$$

For notational simplicity we denote  $\eta_{\text{fmh}} = \eta_{\pi_{\text{FMH}}^*}$ ,  $P_{\text{fmh}} = P_{\pi_{\text{FMH}}^*}$ ,  $\eta_{\text{reg}} = \eta_{\pi_{\text{reg}}^*}$  and  $P_{\text{reg}} = P_{\pi_{\text{reg}}^*}$ . We thus have

$$H = \sum_{s \in \mathcal{S}} \left( \frac{\sigma^2(s)}{\eta_{\text{fmh}}(s)} - \frac{\sigma^2(s)}{\eta_{\text{reg}}(s)} \right) + \rho_n \left( \frac{1}{\gamma(P_{\text{fmh}})} - \frac{1}{\gamma(P_{\text{reg}})} \right) + O(n^{-3/2}).$$

Given that  $\eta_{\text{reg}}$  is a stationary state distribution w.r.t. the MDP dynamics, we can write by optimality of  $\eta^*$  for the problem (7)

$$\sum_{s \in \mathcal{S}} \left( \frac{\sigma^2(s)}{\eta_{\text{fmh}}(s)} - \frac{\sigma^2(s)}{\eta_{\text{reg}}(s)} \right) \leq \sum_{s \in \mathcal{S}} \left( \frac{\sigma^2(s)}{\eta_{\text{fmh}}(s)} - \frac{\sigma^2(s)}{\eta^*(s)} \right) \leq \frac{\sigma_{\max}^2}{\underline{\eta}^2} \|\eta_{\text{fmh}} - \eta^*\|_1 \leq \frac{\sigma_{\max}^2 \sqrt{S}}{\underline{\eta}^2} \|\eta_{\text{fmh}} - \eta^*\|_2.$$

Using successively the triangular inequality, the property guaranteed in  $(\mathcal{P}_2)$  that  $\eta_{\text{fmh}}$  minimizes the distance  $\|\cdot - \eta_1\|_2$  among all the stationary state distributions w.r.t. the MDP dynamics, and finally the property guaranteed in  $(\mathcal{P}_1)$  of  $\delta_n$ -proximity of  $\eta_1$  to  $\eta^*$ , we get

$$\|\eta_{\text{fmh}} - \eta^*\|_2 \leq \|\eta_{\text{fmh}} - \eta_1\|_2 + \|\eta_1 - \eta^*\|_2 \leq 2\|\eta_1 - \eta^*\|_2 \leq 2\delta_n.$$

We conclude the proof using the fact that  $\gamma(P_{\text{fmh}})$  and  $\gamma(P_{\text{reg}})$  are larger than  $\gamma_{\min}$ .  $\square$

## D Proof of Thm. 1

### D.1 Preliminaries

We recall that the notation  $u_n = \tilde{O}(v_n)$  means that there exist  $c > 0$  and  $d > 0$  such that  $u_n \leq c(\log n)^d v_n$  for sufficiently large  $n$ . By abuse of language we say that a stationary policy  $\pi$  belongs to  $\Lambda_\eta$  if  $\forall s \in \mathcal{S}, \eta_\pi \geq 2\underline{\eta}$ . For notational convenience we consider throughout the proof that we relax Asm. 1 (cf. App. A.1) and that the initial state  $s_1$  is drawn from an arbitrary initial distribution over states and we collect its observation  $x_1$ . This leads to the configuration where at every time  $t$  exactly  $t$  state samples have been collected. We start our analysis with the two following technical lemmas.

**Lemma 4.** *Let  $\delta \in (0, 1)$ . For any length  $\tau > 0$ , the following bound holds simultaneously for any state  $s$  and any policy  $\pi \in \Lambda_\eta$  with probability at least  $1 - \delta$*

$$\left| \frac{\sum_{t=1}^{\tau} \mathbb{I}\{\pi_t = s\}}{\tau} - \eta_\pi(s) \right| \leq M(\tau, \delta) := \sqrt{\frac{2B}{\gamma_{\min}\tau}} + \frac{20B}{\gamma_{\min}\tau} \quad \text{with} \quad B = \log\left(\frac{SA^S}{\delta} \sqrt{\frac{1}{\underline{\eta}}}\right).$$

*Proof.* Pick any  $\delta \in (0, 1)$ . Let  $\pi$  be a fixed policy whose stationary distribution is lower-bounded by  $\eta_{\pi, \min}$  and whose associated Markov chain admits  $\gamma_\pi$  as spectral gap. For any length  $\tau > 0$  and state  $s$ , we define  $\nu_{\pi, \tau}(s) = \sum_{t=1}^{\tau} \mathbb{I}\{\pi_t = s\}$ . From Prop. 1, for a fixed state  $s \in \mathcal{S}$ , the following bound holds with probability at least  $1 - \delta$

$$\left| \frac{\nu_{\pi, \tau}(s)}{\tau} - \eta_\pi(s) \right| \leq \sqrt{8\eta_\pi(s)(1 - \eta_\pi(s))\tilde{\epsilon}} + 20\tilde{\epsilon} \quad \text{where} \quad \tilde{\epsilon} = \frac{\log\left(\frac{1}{\delta} \sqrt{\frac{2}{\eta_{\pi, \min}}}\right)}{\gamma_\pi \tau}.$$

Since we need this statement to hold simultaneously for all states  $s \in \mathcal{S}$  and all stationary policies  $\pi \in \Lambda_\eta$ , we need to take a union bound over states and a cover over the action simplex at each state, which leads to tuning  $\delta = \delta'/SA^S$  and thus yields with probability at least  $1 - \delta$

$$\left| \frac{\nu_{\pi, \tau}(s)}{\tau} - \eta_\pi(s) \right| \leq \sqrt{8\eta_\pi(s)(1 - \eta_\pi(s))\tilde{\epsilon}} + 20\tilde{\epsilon} \quad \text{where} \quad \tilde{\epsilon} = \frac{\log\left(\frac{SA^S}{\delta} \sqrt{\frac{1}{\underline{\eta}}}\right)}{\gamma_{\min}\tau}.$$



Using the fact that the function  $x \mapsto x(1-x)$  is upper bounded by  $1/4$  and setting  $B = \log\left(\frac{SA^S}{\delta} \sqrt{\frac{1}{\underline{\eta}}}\right)$  yields the desired high-probability result.  $\square$

**Lemma 5.** *Let  $\delta \in (0, 1)$ . There exists a length  $\tau_\delta > 0$  such that for any  $T \geq \tau_\delta$ , the following inequality holds simultaneously for any state  $s$  and any policy  $\pi \in \Lambda_{\underline{\eta}}$  with probability at least  $1 - \delta$*

$$\sum_{t=1}^T \mathbb{I}\{\pi_t = s\} \geq \underline{\eta}T$$

*Proof.* Pick any  $\delta \in (0, 1)$ .  $M(\tau, \delta)$  is a decreasing function of  $\tau$ , hence there exists a length  $\tau_\delta$  such that for any  $T \geq \tau_\delta$ ,  $M(T, \delta) \leq \underline{\eta}$ . As a result, Lem. 4 guarantees that we have with probability at least  $1 - \delta$  simultaneously for any state  $s$  and any stationary policy  $\pi \in \Lambda_{\underline{\eta}}$

$$\left| \frac{\sum_{t=1}^T \mathbb{I}\{\pi_t = s\}}{T} - \eta_\pi(s) \right| \leq \underline{\eta},$$

which yields in particular

$$\frac{\sum_{t=1}^T \mathbb{I}\{\pi_t = s\}}{T} \geq \eta_\pi(s) - \underline{\eta} \geq \underline{\eta}.$$

$\square$

Restricting our attention to increasing episode lengths in FW-AME and using Lem. 5, we deduce the important property that for any  $\delta \in (0, 1)$ , there exists an episode  $k_\delta$  such that for all episodes  $k$  succeeding it (and including it), we have with probability at least  $1 - \delta$

$$\sum_{a \in \mathcal{A}} \tilde{\lambda}_k(s, a) \geq \underline{\eta}, \quad \forall s \in \mathcal{S}, \quad \forall k \geq k_\delta. \quad (21)$$

More specifically,  $k_\delta$  is the first episode whose length  $\tau_{k_\delta}$  verifies

$$M(\tau_{k_\delta}, \delta) = \sqrt{\frac{2B}{\gamma_{\min} \tau_{k_\delta}}} + \frac{20B}{\gamma_{\min} \tau_{k_\delta}} \leq \underline{\eta} \quad \text{with} \quad B = \log\left(\frac{SA^S}{\delta} \sqrt{\frac{1}{\underline{\eta}}}\right). \quad (22)$$

We proceed by providing time-dependent lower and upper bounds on the true gradient  $\nabla \mathcal{L}$ , which is unknown. We denote by  $\hat{\mathcal{L}}_t^+$  the empirical optimistic approximation of  $\mathcal{L}$  at any time  $t$ , i.e.,

$$\hat{\mathcal{L}}_t^+(\lambda) = \sum_{s \in \mathcal{S}} \frac{1}{\sum_a \lambda(s, a)} \left[ \hat{\sigma}_t^2(s) + 5R^2 \sqrt{\frac{\log\left(\frac{4St}{\delta}\right)}{T_t(s)}} \right] = \hat{\mathcal{L}}_t(\lambda) + \sum_{s \in \mathcal{S}} \frac{\alpha(t, s, \delta)}{\sum_a \lambda(s, a)}.$$

Here we used that  $\nu(s)$  is an observation distribution supported in  $[0, R]$ . We note that this assumption can be easily extended to the general case of sub-Gaussian distributions as done in Carpentier et al. (2011). From Prop. 2, the following inequalities hold with probability at least  $1 - \delta$  for any  $\lambda$ , time  $t$  and state-action pair  $(s, a)$

$$\nabla \hat{\mathcal{L}}_t^+(\lambda)(s, a) = \nabla \hat{\mathcal{L}}_t(\lambda)(s, a) - \frac{\alpha(t, s, \delta)}{(\sum_b \lambda(s, b))^2} \leq \nabla \mathcal{L}(\lambda)(s, a) \leq \nabla \hat{\mathcal{L}}_t(\lambda)(s, a) + \frac{\alpha(t, s, \delta)}{(\sum_b \lambda(s, b))^2}. \quad (23)$$

Finally, let  $T = t_K - 1$  be the final budget (i.e., the time at the end of the final episode  $K - 1$ ). For the sake of clarity and readability, we make the simplification that the logarithmic term  $\log(T)$  behaves as a constant.

## D.2 Core of the proof

We denote by  $\rho_{k+1}$  the approximation error at the end of each episode  $k$  (i.e., at time  $t_{k+1} - 1$ ). Recalling that  $\beta_k = \tau_k / (t_{k+1} - 1)$ , we have

$$\rho_{k+1} = \mathcal{L}(\tilde{\lambda}_{k+1}) - \mathcal{L}(\lambda^*) = \mathcal{L}((1 - \beta_k)\tilde{\lambda}_k + \beta_k\tilde{\psi}_{k+1}) - \mathcal{L}(\lambda^*).$$

Let  $\psi_{k+1}^* = \operatorname{argmin}_{\lambda \in \Lambda_{\underline{\eta}}} \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \lambda \rangle$  be the state-action stationary distribution that “exact” FW would return at episode  $k$ . We have the following series of inequality

$$\begin{aligned} \rho_{k+1} &\leq \mathcal{L}(\tilde{\lambda}_k) - \mathcal{L}(\lambda^*) + \beta_k \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \tilde{\psi}_{k+1} - \tilde{\lambda}_k \rangle + C_{\underline{\eta}} \beta_k^2 \\ &= \mathcal{L}(\tilde{\lambda}_k) - \mathcal{L}(\lambda^*) + \beta_k \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \psi_{k+1}^* - \tilde{\lambda}_k \rangle + C_{\underline{\eta}} \beta_k^2 + \beta_k \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \tilde{\psi}_{k+1} - \psi_{k+1}^* \rangle \\ &\leq \mathcal{L}(\tilde{\lambda}_k) - \mathcal{L}(\lambda^*) + \beta_k \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \lambda^* - \tilde{\lambda}_k \rangle + C_{\underline{\eta}} \beta_k^2 + \beta_k \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \tilde{\psi}_{k+1} - \psi_{k+1}^* \rangle \\ &\leq (1 - \beta_k) \rho_k + C_{\underline{\eta}} \beta_k^2 + \beta_k \underbrace{\langle \nabla \mathcal{L}(\tilde{\lambda}_k), \hat{\psi}_{k+1}^+ - \psi_{k+1}^* \rangle}_{\epsilon_{k+1}} + \beta_k \underbrace{\langle \nabla \mathcal{L}(\tilde{\lambda}_k), \tilde{\psi}_{k+1} - \hat{\psi}_{k+1}^+ \rangle}_{\Delta_{k+1}}, \end{aligned} \quad (24)$$

where the first step follows from the  $C_{\underline{\eta}}$ -smoothness of  $\mathcal{L}$ , the second inequality comes from the FW optimization step and the definition of  $\psi_{k+1}^*$ , which gives  $\langle \nabla \mathcal{L}(\tilde{\lambda}_k), \psi_{k+1}^* - \tilde{\lambda}_k \rangle \leq \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \lambda^* - \tilde{\lambda}_k \rangle$ , the final step follows from the convexity of  $\mathcal{L}$ . The term  $\epsilon_{k+1}$  measures the error due to an inaccurate estimate of the gradient and the term  $\Delta_{k+1}$  refers to the discrepancy between the stationary state-action distribution  $\hat{\psi}_{k+1}^+$  and the empirical frequency  $\tilde{\psi}_{k+1}$  of its realization for  $\tau_k$  steps.

**Step 1 (Bound on error  $\Delta_{k+1}$ ).** For any  $k \geq k_\delta$ , inequality (21) is verified and we can write

$$\langle \nabla \mathcal{L}(\tilde{\lambda}_k), \tilde{\psi}_{k+1} - \hat{\psi}_{k+1}^+ \rangle = \sum_s \frac{-\sigma^2(s)}{(\sum_b \tilde{\lambda}_k(s, b))^2} \sum_a (\tilde{\psi}_{k+1}(s, a) - \hat{\psi}_{k+1}^+(s, a)) \leq \frac{S\sigma_{\max}^2}{\underline{\eta}^2} \left\| \frac{\nu_{k+1}}{\tau_k} - \eta_{\hat{\pi}_{k+1}^+} \right\|_\infty.$$

Let  $B = \log(\frac{SA^S}{\delta} \sqrt{\frac{1}{\underline{\eta}}})$ . From Lem. 4, we have with probability at least  $1 - \delta$  simultaneously for every state  $s$  and every policy followed during the episode

$$\left| \frac{\nu_{k+1}(s)}{\tau_k} - \eta_{\hat{\pi}_{k+1}^+}(s) \right| \leq \sqrt{\frac{2B}{\gamma_{\min} \tau_k}} + \frac{20B}{\gamma_{\min} \tau_k}.$$

Hence we obtain the following bound on  $\Delta_{k+1}$  with probability at least  $1 - \delta$

$$\Delta_{k+1} \leq \frac{S\sigma_{\max}^2}{\underline{\eta}^2} \left[ \sqrt{\frac{2B}{\gamma_{\min} \tau_k}} + \frac{20B}{\gamma_{\min} \tau_k} \right].$$

**Step 2 (Bound on error  $\epsilon_{k+1}$ ).** Using inequality (23), we get with probability at least  $1 - \delta$

$$\begin{aligned} \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \hat{\psi}_{k+1}^+ \rangle &= \sum_{s,a} \hat{\psi}_{k+1}^+(s, a) \nabla \mathcal{L}(\tilde{\lambda}_k)(s, a) \\ &\leq \sum_{s,a} \hat{\psi}_{k+1}^+(s, a) \nabla \hat{\mathcal{L}}_{t_{k-1}}(\tilde{\lambda}_k)(s, a) + \sum_{s,a} \hat{\psi}_{k+1}^+(s, a) \frac{\alpha(t_k - 1, s, \delta)}{(\sum_b \tilde{\lambda}_k(s, b))^2} \\ &\leq \sum_{s,a} \hat{\psi}_{k+1}^+(s, a) \nabla \hat{\mathcal{L}}_{t_{k-1}}^+(\tilde{\lambda}_k)(s, a) + 2 \sum_{s,a} \hat{\psi}_{k+1}^+(s, a) \frac{\alpha(t_k - 1, s, \delta)}{(\sum_b \tilde{\lambda}_k(s, b))^2} \\ &\leq \langle \nabla \hat{\mathcal{L}}_{t_{k-1}}^+(\tilde{\lambda}_k), \psi_{k+1}^* \rangle + 2 \sum_{s,a} \hat{\psi}_{k+1}^+(s, a) \frac{\alpha(t_k - 1, s, \delta)}{(\sum_b \tilde{\lambda}_k(s, b))^2} \\ &\leq \langle \nabla \mathcal{L}(\tilde{\lambda}_k), \psi_{k+1}^* \rangle + 2 \sum_{s,a} \hat{\psi}_{k+1}^+(s, a) \frac{\alpha(t_k - 1, s, \delta)}{(\sum_b \tilde{\lambda}_k(s, b))^2}. \end{aligned}$$

For notational simplicity we denote by  $T_k(s) = T_{t_k-1}(s)$  the number of visits of state  $s$  until the end of episode  $k-1$  (i.e., at time  $t_k-1$ ). Using inequality (21) and an intersection bound over two high-probability events, we get with probability at least  $1-2\delta$  for any episode  $k \geq k_\delta$

$$\begin{aligned} \epsilon_{k+1} &\leq \sum_{s,a} \widehat{\psi}_{k+1}^+(s,a) \frac{10R^2}{\eta^2} \sqrt{\log\left(\frac{4S(t_k-1)}{\delta}\right)} \frac{1}{\sqrt{T_k(s)}} \\ &\leq c_0 \underbrace{\sum_{s,a} \widetilde{\psi}_{k+1}(s,a) \frac{1}{\sqrt{T_k(s)}}}_{v_k} + c_0 \underbrace{\sum_{s,a} (\widehat{\psi}_{k+1}^+(s,a) - \widetilde{\psi}_{k+1}(s,a)) \frac{1}{\sqrt{T_k(s)}}}_{\xi_{k+1}}, \end{aligned}$$

where we define  $c_0 = \frac{10R^2}{\eta^2} \sqrt{\log\left(\frac{4ST}{\delta}\right)}$ .  $\xi_{k+1}$  can be bounded in the same vein as  $\Delta_{k+1}$  using Lem. 4. The error  $\xi_{k+1}$  is of a higher order than  $\Delta_{k+1}$  and for proof simplicity we consider the following loose bound which is satisfied with probability at least  $1-\delta$

$$\xi_{k+1} \leq c_0 \left[ \sqrt{\frac{2B}{\gamma_{\min} \tau_k}} + \frac{20B}{\gamma_{\min} \tau_k} \right].$$

**Step 3 (putting everything together in (24)).** For  $k \geq k_\delta$ , we get with probability at least  $1-2\delta$

$$\Delta_{k+1} + \xi_{k+1} \leq \frac{c_1}{\sqrt{\tau_k}} + \frac{c_2}{\tau_k} \quad \text{with} \quad \begin{cases} c_1 = \left(c_0 + \frac{S\sigma_{\max}^2}{\eta^2}\right) \sqrt{\frac{2B}{\gamma_{\min}}} \\ c_2 = \left(c_0 + \frac{S\sigma_{\max}^2}{\eta^2}\right) \frac{20B}{\gamma_{\min}} \end{cases},$$

which provides the bound for  $k \geq k_\delta$

$$\rho_{k+1} \leq (1-\beta_k)\rho_k + \beta_k \left(\frac{c_1}{\sqrt{\tau_k}} + \frac{c_2}{\tau_k}\right) + C_\eta \beta_k^2 + \beta_k c_0 v_k. \quad (25)$$

Choosing episode lengths satisfying  $t_k = \tau_1(k-1)^3 + 1$  yields

$$\tau_k = t_{k+1} - t_k = \tau_1(3k^2 - 3k + 1) \geq 3\tau_1 k^2 \quad \text{and} \quad \beta_k = \frac{\tau_k}{t_{k+1} - 1} = \frac{3k^2 - 3k + 1}{k^3} \in \left[\frac{1}{k}, \frac{3}{k}\right].$$

Consequently we get

$$\beta_k \left(\frac{c_1}{\sqrt{\tau_k}} + \frac{c_2}{\tau_k}\right) + C_\eta \beta_k^2 \leq \frac{b_\delta}{k^2} \quad \text{with} \quad b_\delta = \frac{\sqrt{3}c_1}{\sqrt{\tau_1}} + \frac{c_2}{\tau_1 k_\delta} + 9C_\eta. \quad (26)$$

Hence the recurrence inequality (25) becomes

$$\rho_{k+1} \leq \left(1 - \frac{1}{k}\right)\rho_k + \frac{b_\delta}{k^2} + \beta_k c_0 v_k. \quad (27)$$

We pick an integer  $q \geq (S/\tau_1)^{1/3} + 1$  such that  $\rho_q \geq 0$  is satisfied.<sup>10</sup> We define the sequence  $(u_n)_{n \geq q}$  as  $u_q = \rho_q$  and

$$u_{n+1} = \left(1 - \frac{1}{n}\right)u_n + \frac{b_\delta}{n^2} + \beta_n c_0 v_n,$$

with  $b_\delta$  the fixed positive constant defined in (26). From inequality (27), we have  $\rho_k \leq u_k$  for  $k \geq k_\delta$  and an immediate induction guarantees the positivity of the sequence  $(u_n)$ . By rearranging we get

$$(n+1)u_{n+1} - nu_n = \frac{-u_n}{n} + \frac{b_\delta(n+1)}{n^2} + (n+1)\beta_n c_0 v_n \leq \frac{b_\delta(n+1)}{n^2} + (n+1)\beta_n c_0 v_n.$$

<sup>10</sup>Assuming this last condition is sensible since as the number of samples increases,  $\tilde{\lambda}$  gets closer to the stationary set  $\Lambda_\eta$  whose minimizer of  $\mathcal{L}$  is  $\lambda^*$ . The introduction of the term  $(S/\tau_1)^{1/3} + 1$  is motivated by the subsequent analysis of the series  $\sum v_k$  in Lem. 6.

By telescoping and using the fact that  $\beta_n \leq 3/n \leq 6/(n+1)$ , we obtain

$$nu_n - qu_q \leq 2b_\delta \sum_{i=q}^{n-1} \frac{1}{i} + 6c_0 \sum_{i=q}^{n-1} v_i \leq 2b_\delta \log\left(\frac{n-1}{q-1}\right) + 6c_0 \sum_{i=q}^{n-1} v_i.$$

Let  $K \geq k_\delta$ . We thus have with probability at least  $1 - 2\delta$

$$\rho_K \leq \frac{q\rho_q + 2b_\delta \log K}{K} + \frac{6c_0}{K} \sum_{k=q}^{K-1} v_k = \frac{\tau_1^{1/3}}{(t_K - 1)^{1/3} + \tau_1^{1/3}} \left( q\rho_q + 2b_\delta \log K + 6c_0 \sum_{k=q}^{K-1} v_k \right). \quad (28)$$

We conclude the proof by plugging the result of Lem. 6 into inequality (28) which yields the desired high-probability bound  $\rho_K = \tilde{O}(1/t_K^{1/3})$ .

**Lemma 6.** *Recalling that  $v_k = \sum_{s,a} \tilde{\psi}_{k+1}(s,a) \frac{1}{\sqrt{T_k(s)}}$ , we have  $\sum v_k = \tilde{O}(1)$ .*

*Proof.* Denoting  $\mathcal{S} = \{1, 2, \dots, S\}$  and recalling that  $q \geq (S/\tau_1)^{1/3} + 1$ , we have

$$\begin{aligned} \sum_{k=q}^{K-1} v_k &= \sum_{k=q}^{K-1} \sum_{s,a} \tilde{\psi}_{k+1}(s,a) \frac{1}{\sqrt{T_k(s)}} = \sum_{k=q}^{K-1} \sum_{s=1}^S \frac{\sqrt{\nu_{k+1}(s)}}{\tau_k} \frac{\sqrt{\nu_{k+1}(s)}}{\sqrt{T_k(s)}} \\ &\leq \sqrt{\sum_{k=q}^{K-1} \sum_{s=1}^S \frac{\nu_{k+1}(s)}{\tau_k^2}} \sqrt{\sum_{k=q}^{K-1} \sum_{s=1}^S \frac{\nu_{k+1}(s)}{T_k(s)}} = \sqrt{\underbrace{\sum_{k=q}^{K-1} \frac{1}{\tau_k}}_{\Sigma_1}} \sqrt{\underbrace{\sum_{k=q}^{K-1} \sum_{s=1}^S \left( \frac{T_{k+1}(s)}{T_k(s)} - 1 \right)}_{\Sigma_2}}, \end{aligned}$$

where the inequality uses the Cauchy-Schwarz inequality on the sum indexed doubly by the episodes and the states. Since the Riemann zeta function of 2 is upper bounded by 3, we have

$$\Sigma_1 \leq \frac{1}{3\tau_1} \sum_{k=q}^{K-1} \frac{1}{k^2} \leq \frac{1}{\tau_1}.$$

There remains to show that  $\Sigma_2 = \tilde{O}(1)$ . We introduce the following related optimization problem. For any  $K \geq q$ , we have  $t_K - 1 \geq S$  since we chose  $q \geq (S/\tau_1)^{1/3} + 1$ . Let  $V^*(K)$  be defined by

$$V^*(K) = \max \sum_{k=q}^{K-1} \sum_{s=1}^S (h_{s,k} - 1), \quad (29)$$

$$\text{s.t. } h_{s,k} \geq 1 \quad \text{and} \quad \sum_{s=1}^S \prod_{k=q}^{K-1} h_{s,k} \leq t_K - 1. \quad (30)$$

We have for any episode  $k$  and state  $s$ ,  $T_{k+1}(s) \geq T_k(s)$  and

$$\sum_{s=1}^S \prod_{k=q}^{K-1} \frac{T_{k+1}(s)}{T_k(s)} = \sum_{s=1}^S \frac{T_K(s)}{T_q(s)} \leq t_K - 1.$$

Hence the sequence  $\left( \frac{T_{k+1}(s)}{T_k(s)} \right)_{s,k}$  satisfies the constraints (30), thus  $\Sigma_2 \leq V^*(K)$ . There remains to solve the optimization problem (29). Since the variables  $h_{s,k}$  play interchangeable roles, there exists  $h^* = h_{s,k}$  for all  $s$  and  $k$ . From the second constraint in (30), we know that  $h^* \leq ((t_K - 1)/S)^{1/(K-q)}$ . Given that (29) is a maximization problem that increases proportionally with  $h^*$ , when  $t_K - 1 \geq S$  (so as to satisfy the first constraint), we finally

have  $h^* = ((t_K - 1)/S)^{1/(K-q)}$ . Consequently we have

$$\Sigma_2 \leq \sum_{k=q}^{K-1} \sum_{s=1}^S \left( \left( \frac{t_K - 1}{S} \right)^{1/(K-q)} - 1 \right) = \underbrace{\frac{\exp\left(\frac{1}{K-q} \log\left(\frac{\tau_1(K-1)^3}{S}\right)\right) - 1}{\frac{1}{K-q} \log\left(\frac{\tau_1(K-1)^3}{S}\right)}}_{\rightarrow 1 \text{ when } K \rightarrow +\infty} \underbrace{S \log\left(\frac{\tau_1(K-1)^3}{S}\right)}_{=\tilde{O}(1)},$$

which proves that  $\Sigma_2 = \tilde{O}(1)$ . We conclude the proof using that  $\sum_{k=q}^{K-1} v_{k+1} \leq \frac{1}{\sqrt{\tau_1}} \sqrt{\Sigma_2}$ . □

---

**Algorithm 3** FW-AME w/ FMH-SDP
 

---

**Input:**  $\tilde{\lambda}_1 = 1/SA, \eta$

**for**  $k = 1, 2, \dots, K - 1$  **do**

$$\hat{\psi}_{k+1}^+ = \operatorname{argmin}_{\lambda \in \Lambda_\eta} \langle \nabla \hat{\mathcal{L}}_{t_k-1}^+(\tilde{\lambda}_k), \lambda \rangle$$

$$\hat{\pi}_{k+1}^+(a|s) = \frac{\hat{\psi}_{k+1}^+(s, a)}{\sum_{b \in \mathcal{A}} \hat{\psi}_{k+1}^+(s, b)}$$

Compute  $\pi_{\text{FMH}}^* = \text{FMH-SDP}(\hat{\pi}_{k+1}^+, \tau_k)$  with  $\delta_{\tau_k}$  defined in Eq. (31)

Execute  $\pi_{\text{FMH}}^*$  for  $\tau_k$  steps, collect the samples and update  $\tilde{\lambda}_{k+1}$  as in Alg. 1

**end for**

---

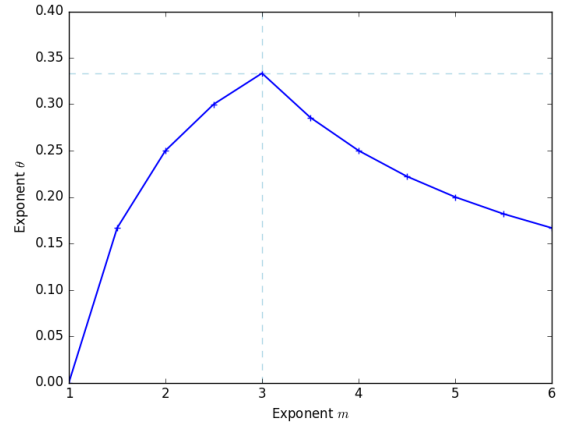


Figure 5: Exponent  $\theta$  as a function of  $m$  (cf. App. D.3).

### D.3 Optimality of the Episode Length

As explained in Sect. 3.2, an interesting open question is whether the regret bound obtained in Thm. 1 is optimal. Our analysis however yields the following optimality result: among all the episode lengths such that the time  $t$  is polynomial in the number of episodes  $k$ , i.e., among all the integers  $m \geq 1$  such that  $t$  behaves as  $k^m$ , the value of  $m$  that optimizes convergence is  $m = 3$ . Indeed, we can apply the Euler method on inequality (25) which results in solving the differential equation  $y' = \frac{-y}{x} + \frac{1}{x^2} + \frac{1}{x^{(m+1)/2}}$  and finding the largest  $\theta$  such that  $x^\theta y(x) = \tilde{O}(1)$ .  $\theta$  is thus the largest value such that

$$\mathcal{L}(\tilde{\lambda}_{k+1}) - \mathcal{L}(\lambda^*) = \tilde{O}\left(\frac{1}{t^\theta}\right) = \tilde{O}\left(\frac{1}{k^{\theta/m}}\right).$$

Fig. 5 plots the exponent  $\theta$  as a function of  $m$  and shows that  $\theta$  reaches its maximal value of  $1/3$  for  $m = 3$ , consequently yielding the regret  $\tilde{O}(1/t^{1/3})$ .

### D.4 FW-AME w/ FMH-SDP

The variant incorporating the framework of FMH-SDP is presented in Alg. 3 and its difference with Alg. 1 is highlighted in blue. The regret analysis is the same as in App. D.2 except that the error  $\Delta_{k+1}$  in the recurrence inequality (24) goes from  $\tilde{\psi}_{k+1} - \hat{\psi}_{k+1}^+$  to

$$\tilde{\psi}_{k+1}^{\text{FMH}} - \hat{\psi}_{k+1}^+ = \tilde{\psi}_{k+1}^{\text{FMH}} - \hat{\psi}_{k+1}^{\text{FMH}} + \hat{\psi}_{k+1}^{\text{FMH}} - \hat{\psi}_{k+1}^+,$$

where  $\widehat{\psi}_{k+1}^{\text{FMH}} = \text{FMH-SDP}(\widehat{\psi}_{k+1}^+, \tau_k)$  and  $\widetilde{\psi}_{k+1}^{\text{FMH}}$  is its empirical realization for the  $\tau_k$  steps of the episode. The new error  $\Delta_{k+1}$  can thus be decomposed as follows

$$\underline{\eta}^2 \Delta_{k+1} \leq \sum_s \sigma^2(s) \left| \frac{\nu_{k+1}^{\text{FMH}}(s)}{\tau_k} - \eta_{\widehat{\pi}_{k+1}^{\text{FMH}}}(s) \right| + \sum_s \sigma^2(s) \left| \eta_{\widehat{\pi}_{k+1}^{\text{FMH}}}(s) - \eta_{\widehat{\pi}_{k+1}^+}(s) \right|,$$

where the first term is  $O\left(1/\sqrt{\gamma(\widehat{\psi}_{k+1}^{\text{FMH}})\tau_k}\right)$  and the second term is upper bounded by  $\sum_s \sigma^2(s)\delta_{\tau_k}$  where  $\delta_{\tau_k}$  is the FMH-SDP parameter from problem (20).

Since w/o FMH-SDP we have  $\gamma_{k+1} = O\left(1/\sqrt{\gamma(\widehat{\psi}_{k+1}^{\text{FMH}})\tau_k}\right)$ , this suggests that the slack variable  $\delta_{\tau_k}$  can decrease at least as  $O(1/\sqrt{\tau_k})$  so as to guarantee that the order of the error  $\Delta_{k+1}$  is unchanged. Furthermore, the component  $\delta_{\tau_k}(s)$  is weighted by  $\sigma^2(s)$  (which is unknown), hence we are encouraged to set

$$\delta_{\tau_k}(s) = \frac{\widehat{\Sigma} - \widehat{\sigma}_{t_{k-1}}^2(s)}{(S-1)\widehat{\Sigma}} \frac{1}{\sqrt{\tau_k}} \quad \text{where } \widehat{\Sigma} = \sum_{s \in \mathcal{S}} \widehat{\sigma}_{t_{k-1}}^2(s). \quad (31)$$

The regret analysis of FW-AME w/ FMH-SDP is thus unchanged and we recover the final rate in  $O(t^{-1/3})$ . In addition, if the heuristic is able to obtain an improvement in the mixing properties of the episodic policy (i.e.,  $\gamma(\widehat{\psi}_{k+1}^{\text{FMH}})$  bigger than  $\gamma(\widehat{\psi}_{k+1}^+)$ ) that outweighs the error introduced by  $\delta_{\tau_k}(s)$ , then the regret performance at episode  $k$  of FW-AME w/ FMH-SDP is improved.

## E Garnet MDPs

We detail here the process for generating Garnet<sup>11</sup> MDPs which we use in Sect. 5. A Garnet instance  $\mathcal{G}(S, A, b, \sigma_{\min}^2, \sigma_{\max}^2)$  is characterized by 5 parameters.  $S$  and  $A$  are the number of states and actions respectively, and  $b$  is a branching factor specifying the number of possible next states for each state-action pair, i.e., the number of uniformly distributed non-zero entries in each line of the MDP transition matrix. We ensure the aperiodicity of the MDP by adding a non-zero probability (equal to 0.001) of self-loop for all state-action pairs. Since the state means are arbitrarily fixed, there remains to uniformly sample the state variances  $\sigma^2(s)$  between  $\sigma_{\min}^2$  and  $\sigma_{\max}^2$  and randomly select two states whose variances are set respectively to  $\sigma_{\min}^2$  and  $\sigma_{\max}^2$ . We likewise introduce reversible Garnet MDPs denoted by  $\mathcal{G}_{\mathcal{R}}$ . The generation process of  $\mathcal{G}_{\mathcal{R}}$  is identical to  $\mathcal{G}$  except that we set the branching factor to  $b-1$  and ensure the reversibility of the MDP by randomly picking  $a \in \mathcal{A}$  and  $q \in (0, 1)$  such that  $p(s|s', a) = q$  for every pair  $(s, s')$  such that  $Q(s, s') = 1$  (and finally normalize to obtain an admissible  $p$ ).

We note that the Garnet procedure allows some control over the mixing properties of the MDP. Indeed, when  $A$  and  $b$  are small, only a few transitions are assigned significant probabilities so the speed of mixing is generally slower. For higher values of  $A$  and  $b$ , all the positive transition probabilities are of similar magnitude so the speed of mixing is generally faster.

<sup>11</sup>In full, Generalized Average Reward Non-stationary Environment Test-bench (Bhatnagar et al., 2009).