A Previous Work on Solving DMDP with a Full Model

Value iteration was proposed by [Bel57] to compute an exact optimal policy of a given DMDP in time $\mathcal{O}((1-\gamma)^{-1}|\mathcal{S}|^2|\mathcal{A}|L\log((1-\gamma)^{-1}))$, where L is the total number of bits needed to represent the input; and it can find an approximate ϵ -approximate solution in time $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|(1 - \epsilon))$ γ)⁻¹log(1/ ϵ (1 - γ))); see e.g. [Tse90, LDK95]. The policy iteration was introduced by [How60] shortly after, where the policy is monotonically improved according to its associated value function. Its complexity has also been analyzed extensively; see e.g. [MS99, Ye11, Sch13]. Ye [Ye11] showed that policy iteration and the simplex method are strongly polynomial for DMDP and terminates in $\mathcal{O}(|\mathcal{S}|^2|\mathcal{A}|(1-\gamma)^{-1}\log(|\mathcal{S}|(1-\gamma)^{-1}))$ number of iterations. Later [HMZ13] and [Sch13] improved the iteration bound to $O(|\mathcal{S}||\mathcal{A}|(1-\gamma)^{-1}\log((1-\gamma)^{1}))$ for Howard's policy iteration method. A third approach is to formulate the nonlinear Bellman equation into a linear program [d'E63, DG60], and solve it using standard linear program solvers, such as the simplex method by Dantzig [Dan16] and the combinatorial interior-point algorithm by [Ye05]. [LS14, LS15] showed that one can solve linear programs in $O(\sqrt{\operatorname{rank}(A)})$ number of linear system solves, which, applied to DMDP, yields to a running time of $\widetilde{O}(|\mathcal{S}|^{2.5}|\mathcal{A}|L)$ for computing the exact policy and $\widetilde{O}(|\mathcal{S}|^{2.5}|\mathcal{A}|\log(1/\epsilon))$ for computing an ϵ -optimal policy. [SWWY18] further improved the complexity of value iteration by using randomization and variance reduction. See Table 2 for comparable run-time results or computing the optimal policy when the MDP model is fully given.

Algorithm	Complexity	References
Value Iteration (exact)	$ \mathcal{S} ^2 \mathcal{A} L rac{\log(1/(1-\gamma))}{1-\gamma}$	[Tse90, LDK95]
Value Iteration	$ \mathcal{S} ^2 \mathcal{A} rac{\log(1/(1-\gamma)\epsilon)}{1-\gamma}$	[Tse90, LDK95]
Policy Iteration (Block Simplex)	$\frac{ \mathcal{S} ^4 \mathcal{A} ^2}{1-\gamma}\log\bigl(\frac{1}{1-\gamma}\bigr)$	[Ye11],[Sch13]
Recent Interior Point Methods	$\widetilde{O}(\mathcal{S} ^{2.5} \mathcal{A} L) \ \widetilde{O}(\mathcal{S} ^{2.5} \mathcal{A} \log(1/\epsilon))$	[LS14]
Combinatorial Interior Point Algorithm	$ \mathcal{S} ^4 \mathcal{A} ^4 \log rac{ \mathcal{S} }{1-\gamma}$	[Ye05]
High Precision Randomized Value Iteration	$\widetilde{O}\left[\left(\operatorname{nnz}(P) + \frac{ \mathcal{S} \mathcal{A} }{(1-\gamma)^3}\right)\log\left(\frac{1}{\epsilon\delta}\right)\right]$	[SWWY18]

Table 2: Running Times to Solve DMDPs Given the Full MDP Model: In this table, |S| is the number of states, |A| is the number of actions per state, $\gamma \in (0, 1)$ is the discount factor, and L is a complexity measure of the linear program formulation that is at most the total bit size to present the DMDP input. Rewards are bounded between 0 and 1.

B Sample and Time Efficient Value Computation

In this section, we describe an algorithm that obtains an ϵ -optimal values in time $\widetilde{O}(\epsilon^{-2}(1 - \gamma)^{-3}|\mathcal{S}||\mathcal{A}|)$. Note that the time and number of samples of this algorithm is optimal (up to logarithmic factors) due to the lower bound in [AMK13] which also established this upper bound on the sample complexity (but not time complexity) of the problem.

We achieve this by combining the algorithms in [AMK13] and [SWWY18]. First, we use the ideas and analysis of [AMK13] to construct a sparse MDP where the optimal value function of this MDP approximates the optimal value function of the original MDP and then we run the high precision algorithm in [SWWY18] on this sparsified MDP. We show that [SWWY18] runs in nearly linear time on sparsified MDP. Since the number of samples taken to construct the sparsified MDP was the the optimal number of samples, to solve the problem, the ultimate running time we thereby achieve is nearly optimal as any algorithm needs spend time at least the number of samples to obtain these samples.

We include this for completeness but note that the approximate value function we show how to compute here does not suffice to compute policy of the MDP of comparable quality. The greedy policy of an ϵ -optimal value function is an $\epsilon/(1 - \gamma)$ -optimal policy in the worst case. It has been shown in [AMK13] that the greedy policy of their value function is ϵ -optimal if $\epsilon \leq (1 - \gamma)^{1/2} |S|^{-1/2}$. However, when ϵ is so small, the seemingly sublinear runtime $\widetilde{O}((1 - \gamma)^{-3}S||\mathcal{A}|/\epsilon^2)$ essentially means a linear running time and sample complexity as $O((1 - \gamma)^{-3}|\mathcal{S}|^2|\mathcal{A}|)$. The running time can be obtained by merely applying the result in [SWWY18] (although with a slightly different computation model).

B.1 The Sparsified DMDP

Suppose we are given a DMDP $\mathcal{M} = (S, \mathcal{A}, \mathbf{r}, \mathbf{P}, \gamma)$ with a sampling oracle. To approximate the optimal value of this MDP, we perform a spasification procedure as in [AMK13]. Sparsification of DMDP is conducted as follows. Let $\delta > 0, \epsilon > 0$ be arbitrary. First we pick a number

$$m = \Theta\left[\frac{1}{(1-\gamma)^{3}\epsilon^{2}}\log\left(\frac{|\mathcal{S}||\mathcal{A}|}{\delta}\right)\right] .$$
(B.1)

For each $s \in S$ and each $a \in A$, we generate a sequence of independent samples from S using the probability vector $P_{s,a}$

$$s_{s,a}^{(1)}, s_{s,a}^{(2)}, \dots, s_{s,a}^{(m)}$$

Next we construct a new and sparse probability vector $\widehat{P}_{s,a} \in \Delta_{|S|}$ as

$$\forall s' \in \mathcal{S} : \widehat{\boldsymbol{P}}_{s,a}(s') = \frac{1}{m} \cdot \sum_{i=1}^{m} \mathbf{1}(s_{s,a}^{(i)} = s').$$

Combining these |S||A| new probability vectors, we obtain a new probability transition matrix $\hat{P} \in \mathbb{R}^{S \times A \times S}$ with number of non-zeros

$$\operatorname{nnz}(\widehat{\boldsymbol{P}}) = O\left[\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3 \epsilon^2} \log\left(\frac{|\mathcal{S}||\mathcal{A}|}{\delta}\right)\right]$$

Denote $\widehat{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, \mathbf{r}, \widehat{\mathbf{P}}, \gamma)$ as the sparsified DMDP. In the rest of this section, we use $\widehat{\cdot}$ to represent the quantities corresponding to DMDP $\widehat{\mathcal{M}}$, e.g., $\widehat{\boldsymbol{v}}^*$ for the optimal value function, $\widehat{\pi}^*$ for a optimal policy, and $\widehat{\boldsymbol{Q}}^*$ for the optimal *Q*-function. There is a strong approximation guarantee of the optimal *Q*-function of the sparsified MDP, presented as follows.

Theorem B.1 ([AMK13]). Let \mathcal{M} be the original DMDP and $\widehat{\mathcal{M}}$ be the corresponding sparsified version. Let Q^* be the optimal Q-function vector of the original DMDP and \widehat{Q}^* be the optimal Q-function of $\widehat{\mathcal{M}}$. Then with probability at least $1 - \delta$ (over the randomness of the samples),

$$\|\widehat{\boldsymbol{Q}}^* - \boldsymbol{Q}^*\|_{\infty} \leq \epsilon$$

Recall that v^* and \hat{v}^* are the optimal value functions of \mathcal{M} and $\widehat{\mathcal{M}}$. From Theorem B.1, we immediately have

$$\forall s \in \mathcal{S} : |\boldsymbol{v}^*(s) - \widehat{\boldsymbol{v}}^*(s)| = |\max_{a \in \mathcal{A}} \boldsymbol{Q}^*(s, a) - \max_{a \in \mathcal{A}} \widehat{\boldsymbol{Q}}^*(s, a)| \le \max_{a \in \mathcal{A}} |\boldsymbol{Q}^*(s, a) - \widehat{\boldsymbol{Q}}^*(s, a)| \le \epsilon,$$

with probability at least $1 - \delta$.

B.2 High Precision Algorithm in the Sparsified MDP

Next we shall use the high precision algorithm of the [SWWY18] which has the following guarantee. **Theorem B.2** ([SWWY18]). *There is an algorithm which given an input DMDP* $\mathcal{M} = (S, \mathcal{A}, \mathbf{r}, \mathbf{P}, \gamma)$ in time⁶

$$\widetilde{O}\left[\left(\operatorname{nnz}(\boldsymbol{P}) + \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3}\right) \cdot \log \epsilon^{-1} \cdot \log \delta^{-1}\right]$$

and outputs a vector \tilde{v}^* such that with probability at least $1 - \delta$,

$$\|\widetilde{\boldsymbol{v}}^* - \boldsymbol{v}^*\|_{\infty} \leq \epsilon.$$

where v^* is the optimal value of \mathcal{M} .

 $^{6}\widetilde{O}(f)$ denotes $O(f \cdot \log^{O(1)} f)$.

Combining the above two theorems, we immediately obtain an algorithm for finding ϵ -optimal value functions. It works by first generating enough samples for each state-action pair and then call the high-precision MDP solver by [SWWY18]. It does not sample transitions adaptively. We show that it achieves an optimal running time guarantee (up to poly log factors) of obtaining the value function under the sampling oracle model.

Theorem B.3. Given an input DMDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbf{r}, \mathbf{P}, \gamma)$ with a sampling oracle and optimal value function \mathbf{v}^* , there exists an algorithm, that runs in time

$$\widetilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3} \cdot \frac{1}{\epsilon^2} \cdot \log^2\left(\frac{1}{\delta}\right)\right)$$

and outputs a vector \hat{v}^* such that $\|\hat{v}^* - v^*\|_{\infty} \leq O(\epsilon)$ with probability at least $1 - O(\delta)$.

Proof. We first obtain a sparsified MDP $\widehat{\mathcal{M}} = (\mathcal{S}, \mathcal{A}, \mathbf{r}, \widehat{\mathbf{P}}, \gamma)$ using the procedure described in Section B.1. This procedure runs in time $O(|\mathcal{S}||\mathcal{A}|m)$, recalling that m is the number of samples per (s, a), defined in (B.1). Let $\widehat{\mathbf{u}}^*$ be the optimal value function of $\widehat{\mathcal{M}}$. By Theorem B.1, with probability at least $1 - \delta$, $\|\widehat{\mathbf{u}}^* - \mathbf{v}^*\| \le \epsilon$, which we condition on for the rest of the proof. Calling the algorithm in Theorem B.2, we obtain a vector $\widetilde{\mathbf{u}}^*$ in time

$$\widetilde{O}\left[\left(\operatorname{nnz}(\widehat{\boldsymbol{P}}) + \frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3}\right) \cdot \log \epsilon^{-1} \cdot \log \delta^{-1}\right] = \widetilde{O}\left(\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3} \cdot \frac{1}{\epsilon^2} \cdot \log^2 \frac{1}{\delta}\right)$$

and that with probability at least $1-\delta$, $\|\widetilde{\boldsymbol{u}}^* - \widehat{\boldsymbol{u}}^*\| \le \epsilon$, which we condition on. By triangle inequality, we have

$$\|\widetilde{\boldsymbol{u}}^* - \boldsymbol{v}^*\|_{\infty} \leq \|\widetilde{\boldsymbol{u}}^* - \widehat{\boldsymbol{u}}^*\|_{\infty} + \|\widehat{\boldsymbol{u}}^* - \boldsymbol{v}^*\|_{\infty} \leq 2\epsilon.$$

This concludes the proof.

C Variance Bounds

In this section, we study some properties of a DMDP. Most of the content in this section is similar to [AMK13]. We provide slight modifications and improvement to make the results fit to our application. The main result of this section is to show the following lemma.

Lemma C.1 (Upper Bound on Variance). For any π , we have

$$\|(\boldsymbol{I}-\gamma \boldsymbol{P}^{\pi})^{-1}\sqrt{\boldsymbol{\sigma}_{\boldsymbol{v}^{\pi}}}\|_{\infty}^{2} \leq \frac{1+\gamma}{\gamma^{2}(1-\gamma)^{3}},$$

where $\sigma_{v^{\pi}} = \mathbf{P}^{\pi}(\mathbf{v}^{\pi})^2 - (\mathbf{P}^{\pi}\mathbf{v}^{\pi})^2$ is the "one-step" variance of playing policy π .

Before we prove this lemma, we introduce another notation. We define $\Sigma^{\pi} \in \mathbb{R}^{|S||A|}$ for all $(s, a) \in S \times A$ by

$$\boldsymbol{\Sigma}^{\pi}(s,a) := \mathbb{E}\left[\left(\boldsymbol{r}(s,a) + \sum_{t \ge 1} \gamma^{t} \boldsymbol{r}(s^{t},a^{t}) - \boldsymbol{Q}^{\pi}(s,a)\right)^{2} \middle| s^{0} = s, a^{0} = a, a^{t} = \pi(s^{t})\right]$$

where $a^t = \pi(s^t)$. Thus Σ^{π} is the variance of the reward of starting with (s, a) and play π for infinite steps. The crucial observation of obtaining the near-optimal sample complexity is the following "Bellman Equation" for variance. It is a consequence of "the law of total variance".

Lemma C.2 (Bellman Equation for variance). Σ^{π} satisfies the Bellman equation

$$oldsymbol{\Sigma}^{\pi} = \gamma^2 oldsymbol{\sigma}_{oldsymbol{v}^{\pi}} + \gamma^2 \cdot oldsymbol{P}^{\pi} oldsymbol{\Sigma}^{\pi}.$$

Proof. By direct expansion,

$$\boldsymbol{\Sigma}^{\pi}(s,a) = \mathbb{E}\left[\left(\boldsymbol{r}(s,a) + \sum_{t \ge 1} \gamma^{t} \boldsymbol{r}(s^{t},a^{t})\right)^{2} \middle| s^{0} = s, a^{0} = a, a^{t} = \pi(s^{t})\right] - (\boldsymbol{Q}^{\pi}(s,a))^{2}.$$
(C.1)

The first term in RHS can be written as

$$\begin{split} & \mathbb{E}\Big[\left(\boldsymbol{r}(s,a) + \sum_{t \ge 1} \gamma^{t} \boldsymbol{r}(s^{t},a^{t})\right)^{2} \Big| s^{0} = s, a^{0} = a, a^{t} = \pi(s^{t})\Big] \\ &= \sum_{s' \in \mathcal{S}} \boldsymbol{P}_{s,a}(s') \mathbb{E}\Big[\left(\boldsymbol{r}(s,a) + \gamma \boldsymbol{r}(s',\pi(s')) + \gamma \sum_{t \ge 1} \gamma^{t} \boldsymbol{r}(s^{t},a^{t})\right)^{2} \Big| s^{0} = s', a^{0} = \pi(s'), a^{t} = \pi(s^{t})\Big] \\ &= \boldsymbol{r}(s,a)^{2} + 2\gamma \boldsymbol{r}(s,a) \cdot \sum_{s' \in \mathcal{S}} \boldsymbol{P}_{s,a}(s') \boldsymbol{Q}^{\pi}(s',\pi(s')) \\ &\quad + \gamma^{2} \sum_{s' \in \mathcal{S}} \boldsymbol{P}_{s,a}(s') \mathbb{E}\Big[\left(\boldsymbol{r}(s',\pi(s')) + \sum_{t \ge 1} \gamma^{t} \boldsymbol{r}(s^{t},a^{t})\right)^{2} \Big| s^{0} = s', a^{0} = \pi(s'), a^{t} = \pi(s^{t})\Big] \\ &= \boldsymbol{r}(s,a)^{2} + 2\gamma \boldsymbol{r}(s,a) \cdot \sum_{s' \in \mathcal{S}} \boldsymbol{P}_{s,a}(s') \boldsymbol{Q}^{\pi}(s',\pi(s')) + \gamma^{2} (\boldsymbol{P}^{\pi} \boldsymbol{\Sigma}^{\pi})(s,a) + \gamma^{2} \sum_{s' \in \mathcal{S}} \boldsymbol{P}_{s,a}(s') (\boldsymbol{Q}^{\pi}(s',\pi(s')))^{2} \\ &= \boldsymbol{Q}^{\pi}(s,a)^{2} + \gamma^{2} (\boldsymbol{P}^{\pi} \boldsymbol{\Sigma}^{\pi})(s,a) + \gamma^{2} \sum_{s' \in \mathcal{S}} \boldsymbol{P}_{s,a}(s') (\boldsymbol{Q}^{\pi}(s',\pi(s')))^{2} - \gamma^{2} \left(\sum_{s' \in \mathcal{S}} \boldsymbol{P}_{s,a}(s') \boldsymbol{Q}^{\pi}(s',\pi(s'))\right)^{2} \\ &= \boldsymbol{Q}^{\pi}(s,a)^{2} + \gamma^{2} (\boldsymbol{P}^{\pi} \boldsymbol{\Sigma}^{\pi})(s,a) + \gamma^{2} \boldsymbol{\sigma}_{\boldsymbol{v}^{\pi}}(s,a). \end{split}$$

Combining the above two equations, we conclude the proof.

As a remark, we note that

$$\boldsymbol{\Sigma}^{\pi} = \gamma^2 (\boldsymbol{I} - \gamma^2 \boldsymbol{P}^{\pi})^{-1} \boldsymbol{\sigma}_{\boldsymbol{v}^{\pi}}.$$

Furthermore, by definition, we have

$$\max_{(s,a)\in\mathcal{S}} \Sigma^{\pi}(s,a) \le (1-\gamma)^{-2},$$

The next lemma is crucial in proving the error bounds.

Lemma C.3. Let $P \in \mathbb{R}^{n \times n}$ be a non-negative matrix in which every row has ℓ_1 norm at most 1, *i.e.* ℓ_{∞} operator norm at most 1. Then for all $\gamma \in (0, 1)$ and $v \in \mathbb{R}^n_{\geq 0}$ we have

$$\|(\boldsymbol{I}-\gamma\boldsymbol{P})^{-1}\sqrt{\boldsymbol{v}}\|_{\infty} \leq \sqrt{\frac{1}{1-\gamma}\left\|(\boldsymbol{I}-\gamma\boldsymbol{P})^{-1}\boldsymbol{v}\right\|_{\infty}} \leq \sqrt{\frac{1+\gamma}{1-\gamma}\left\|(\boldsymbol{I}-\gamma^{2}\boldsymbol{P})^{-1}\boldsymbol{v}\right\|_{\infty}}$$

Proof. Since, every row of P has ℓ_1 norm at most 1, by Cauchy-Schwarz for $i \in [n]$ we have

$$[oldsymbol{P}\sqrt{oldsymbol{v}}]_i = \sum_{j\in[n]}oldsymbol{P}_{ij}\sqrt{oldsymbol{v}}_j \leq \sqrt{\sum_{j\in[n]}oldsymbol{P}_{ij}\cdot\sum_{j\in[n]}oldsymbol{P}_{ij}oldsymbol{v}_j} \leq \sqrt{oldsymbol{P}oldsymbol{v}} \;.$$

Since v is non-negative and applying P preserves non-negativity, applying this inequality repeatedly yields that $P^k \sqrt{v} \le \sqrt{P^k v}$ entrywise for all k > 0. Consequently, Cauchy-Schwarz again yields

$$(\boldsymbol{I} - \gamma \boldsymbol{P})^{-1} \sqrt{\boldsymbol{v}} = \sum_{i=0}^{\infty} [\gamma \boldsymbol{P}]^i \sqrt{\boldsymbol{v}} \le \sum_{i=0}^{\infty} \gamma^i \sqrt{\boldsymbol{P}^i \boldsymbol{v}} \le \sqrt{\sum_{i=0}^{\infty} \gamma^i \cdot \sum_{i=0}^{\infty} \gamma^i \boldsymbol{P}^i \boldsymbol{v}} \le \sqrt{\frac{1}{1 - \gamma} \| (\boldsymbol{I} - \gamma \boldsymbol{P})^{-1} \boldsymbol{v} \|_{\infty}} \,.$$

Next, as $(I - \gamma P)(I + \gamma P) = (I - \gamma P^2)$ we see that $(I - \gamma P)^{-1} = (I + \gamma P)(I - \gamma^2 P)^{-1}$. Furthermore, as $\|Px\|_{\infty} \leq \|x\|_{\infty}$ for all x we have $\|(I + \gamma P)x\|_{\infty} \leq (1 + \gamma)\|x\|_{\infty}$ for all x and therefore $\|(I - \gamma P)^{-1}v\|_{\infty} \leq (1 + \gamma)\|(I - \gamma^2 P)^{-1}v\|_{\infty}$ as desired. \Box

We are now ready to prove Lemma C.1.

Proof of Lemma C.1. The lemma follows directly from the application of Lemma C.3. This proof is slightly simpler, tighter, and more general than the one in [AMK13]. \Box

D Lower Bounds on Policy

Lemma D.1. Suppose $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, \gamma, \mathbf{r})$ is a DMDP with an sampling oracle. Suppose π is a given policy. Then there is an algorithm, halts in $\widetilde{O}((1-\gamma)^{-3}\epsilon^{-2}|\mathcal{S}|)$ time, outputs a vector \mathbf{v} such that, with high probability, $\|\mathbf{v}^{\pi} - \mathbf{v}\|_{\infty} \leq \epsilon$.

Proof. The lemma follows from a direct application of Theorem B.2.

Remark D.2. Suppose $|\mathcal{A}| = \widetilde{\Omega}(1)$. Suppose there is an algorithm that obtains an ϵ -optimal policy with Z samples, then the above lemma implies an algorithm for obtaining an ϵ -optimal value function with $Z + \widetilde{O}((1 - \gamma)^{-3}\epsilon^{-2}|\mathcal{S}|)$ samples. By the $\Omega((1 - \gamma)^{-3}\epsilon^{-2}|\mathcal{S}||\mathcal{A}|)$ sample bound on obtaining approximate value functions given in [AMK13], the above lemma implies a

 $Z = \Omega((1-\gamma)^{-3}\epsilon^{-2}|\mathcal{S}||\mathcal{A}|) - \widetilde{O}((1-\gamma)^{-3}\epsilon^{-2}|\mathcal{S}|) = \Omega((1-\gamma)^{-3}\epsilon^{-2}|\mathcal{S}||\mathcal{A}|)$

sample lower bound for obtaining an ϵ -optimal policy.

E Missing Proofs

Here are several standard properties of the Bellman value operator (see, e.g., [Ber13]).

Fact 1. Let $v_1, v_2 \in \mathbb{R}^S$ be two vectors. Let \mathcal{T} be a value operator of a DMDP with discount factor γ . Let $\pi \in \mathcal{A}^S$ be an arbitrary policy. Then the follows hold.

- Monotonicity: If $v_1 \leq v_2$ then $\mathcal{T}(v_1) \leq \mathcal{T}(v_2)$;
- Contraction: $\|\mathcal{T}(v_1) \mathcal{T}(v_2)\|_{\infty} \leq \gamma \|v_1 v_2\|_{\infty}$ and $\|\mathcal{T}_{\pi}(v_1) \mathcal{T}_{\pi}(v_2)\|_{\infty} \leq \gamma \|v_1 v_2\|_{\infty}$.

E.1 Missing Proofs from Section 4

To begin, we introduce two standard concentration results. Let $p \in \Delta_{\mathcal{S}}$ be a probability vector, and $v \in \mathbb{R}^{\mathcal{S}}$ be a vector. Let $p_m \in \Delta_{\mathcal{S}}$ be empirical estimations of p using m i.i.d. samples from the distribution p. For instance, let these samples be $s_1, s_2, \ldots, s_m \in \mathcal{S}$, then $\forall s \in \mathcal{S} : p_m(s) = \sum_{j=1}^m \mathbf{1}(s_j = s)/m$.

Theorem E.1 (Hoeffding Inequality). Let $\delta \in (0, 1)$ be a parameter, vectors $\boldsymbol{p}, \boldsymbol{p}_m$ and \boldsymbol{v} defined above. Then with probability at least $1 - \delta$,

$$\left| \boldsymbol{p}^{\top} \boldsymbol{v} - \boldsymbol{p}_m^{\top} \boldsymbol{v} \right| \leq \| \boldsymbol{v} \|_{\infty} \cdot \sqrt{2m^{-1} \log(2\delta^{-1})}.$$

Theorem E.2 (Bernstein Inequality). Let $\delta \in (0, 1)$ be a parameter, vectors $\boldsymbol{p}, \boldsymbol{p}_m$ and \boldsymbol{v} defined as in Theorem E.1. Then with probability at least $1 - \delta$

$$\left|\boldsymbol{p}^{\top}\boldsymbol{v}-\boldsymbol{p}_{m}^{\top}\boldsymbol{v}\right| \leq \sqrt{2m^{-1}\operatorname{Var}_{s'\sim\boldsymbol{p}}(\boldsymbol{v}(s'))\cdot\log(2\delta^{-1})} + (2/3)m^{-1}\|\boldsymbol{v}\|_{\infty}\cdot\log(2\delta^{-1}),$$

where $\operatorname{Var}_{s' \sim p}(\boldsymbol{v}(s')) = \boldsymbol{p}^{\top} \boldsymbol{v}^2 - (\boldsymbol{p}^{\top} \boldsymbol{v})^2$.

Proof of Lemma 4.1. By Theorem E.2 and a union bound over all (s, a) pairs, with probability at least $1 - \delta/4$, for every (s, a), we have

$$\left|\widetilde{\boldsymbol{w}}(s,a) - \boldsymbol{P}_{s,a}^{\top} \boldsymbol{v}^{(0)}\right| \leq \sqrt{2\sigma_{\boldsymbol{v}^{(0)}} \cdot m_1^{-1} \cdot L} + 2 \cdot (3m_1)^{-1} \cdot \|\boldsymbol{v}^{(0)}\|_{\infty} \cdot L, \qquad (E.1)$$

which is the first inequality.

Next, by Theorem E.1 and a union bound over all (s, a) pairs, with probability at least $1 - \delta/4$, for every (s, a), we have

$$\left|\widetilde{\boldsymbol{w}}(s,a) - \boldsymbol{P}_{s,a}^{\top} \boldsymbol{v}^{(0)}\right| \leq \|\boldsymbol{v}^{(0)}\|_{\infty} \cdot \sqrt{2m_1^{-1}L},$$

which we condition on. Thus

$$\begin{split} \left| \widetilde{\boldsymbol{w}}(s,a)^{2} - (\boldsymbol{P}_{s,a}^{\top}\boldsymbol{v}^{(0)})^{2} \right| &= (\widetilde{\boldsymbol{w}}(s,a) + \boldsymbol{P}_{s,a}^{\top}\boldsymbol{v}^{(0)}) \cdot \left| \widetilde{\boldsymbol{w}}(s,a) - \boldsymbol{P}_{s,a}^{\top}\boldsymbol{v}^{(0)} \right| \\ &\leq \left[2\boldsymbol{P}_{s,a}^{\top}\boldsymbol{v}^{(0)} + \|\boldsymbol{v}^{(0)}\|_{\infty} \cdot \sqrt{2m_{1}^{-1}L} \right] \cdot \left| \widetilde{\boldsymbol{w}}(s,a) - \boldsymbol{P}_{s,a}^{\top}\boldsymbol{v}^{(0)} \right| \\ &\leq 2(\boldsymbol{P}_{s,a}^{\top}\boldsymbol{v}^{(0)}) \cdot \|\boldsymbol{v}^{(0)}\|_{\infty} \cdot \sqrt{2m_{1}^{-1}L} + \|\boldsymbol{v}^{(0)}\|_{\infty}^{2} \cdot 2m_{1}^{-1}L. \end{split}$$

Since $\boldsymbol{P}_{s,a}^{\top} \boldsymbol{v}^{(0)} \leq \| \boldsymbol{v}^{(0)} \|_{\infty}$, we obtain

$$\left|\widetilde{\boldsymbol{w}}(s,a)^2 - (\boldsymbol{P}_{s,a}^{\top}\boldsymbol{v}^{(0)})^2\right| \leq 3 \|\boldsymbol{v}^{(0)}\|_{\infty}^2 \cdot \sqrt{2m_1^{-1}L},$$

provided $2m_1^{-1}L \leq 1$. Next by Lemma E.1 and a union bound over all (s, a) pairs, with probability at least $1 - \delta/4$, for every (s, a), we have

$$\left. \frac{1}{m_1} \sum_{j=1}^{m_1} \boldsymbol{v}^2(s_{s,a}^{(j)}) - \boldsymbol{P}_{s,a}^\top \boldsymbol{v}^2 \right| \le \|\boldsymbol{v}^{(0)}\|_{\infty}^2 \cdot \sqrt{2L/m_1}.$$

By a union bound, we obtain, with probability at least $1 - \delta/2$,

$$\begin{aligned} \left| \widehat{\boldsymbol{\sigma}}(s,a) - \boldsymbol{\sigma}_{\boldsymbol{v}^{(0)}}(s,a) \right| &\leq \left| \widetilde{\boldsymbol{w}}(s,a)^2 - (\boldsymbol{P}_{s,a}^{\top} \boldsymbol{v}^{(0)})^2 \right| + \left| m_1^{-1} \sum_{j=1}^{m_1} \boldsymbol{v}^2(s_{s,a}^{(j)}) - \boldsymbol{P}_{s,a}^{\top} \boldsymbol{v}^2 \right| \\ &\leq 4 \| \boldsymbol{v}^{(0)} \|_{\infty}^2 \cdot \sqrt{2m_1^{-1}L}. \end{aligned} \tag{E.2}$$

By a union bound, with probability at least $1-\delta$, both (E.1) and (E.2) hold, concluding the proof.

Proof of Lemma 4.2. Since for each (s, a), $\sigma_{v}(s, a)$ is a variance, then we have triangle inequality, $\sqrt{\sigma_{v}} \leq \sqrt{\sigma_{v^*}} + \sqrt{\sigma_{v-v^*}}$.

Observing that

$$\mathbf{v} - \mathbf{v}^*(s, a) \leq \mathbf{P}_{s, a}^{\top} (\mathbf{v} - \mathbf{v}^*)^2 \leq \epsilon^2 \cdot \mathbf{1}.$$

We conclude the proof by taking a square root of all three sides of the above inequality. \Box

Proof of Lemma 4.3. Recall that for each $(s, a) \in S \times A$,

 σ

$$\boldsymbol{g}^{(i)}(s,a) = \frac{1}{m_2} \sum_{j=1}^{m_2} \left[\boldsymbol{v}^{(i)}(s_{s,a}^{(j)}) - \boldsymbol{v}^{(0)}(s_{s,a}^{(j)}) \right] - (1-\gamma) \frac{u}{8} ,$$

where $m_2 = 128(1-\gamma)^{-2} \cdot \log(2|\mathcal{S}||\mathcal{A}|R/\delta)$ and $s_{s,a}^{(1)}, s_{s,a}^{(2)}, \ldots, s_{s,a}^{(m_2)}$ is a sequence of independent samples from $P_{s,a}$. Thus by Theorem E.1 and a union bound over $\mathcal{S} \times \mathcal{A}$, with probability at least $1-\delta/R$, we have

$$\begin{aligned} \forall (s,a) \in \mathcal{S} \times \mathcal{A} : \left| \sum_{j=1}^{m_2} \left[\boldsymbol{v}^{(i)}(s_{s,a}^{(j)}) - \boldsymbol{v}^{(0)}(s_{s,a}^{(j)}) \right] - \boldsymbol{P}_{s,a}^{\top} \left[\boldsymbol{v}^{(i)} - \boldsymbol{v}^{(0)} \right] \right| \\ & \leq \| \boldsymbol{v}^{(i)} - \boldsymbol{v}^{(0)} \|_{\infty} \sqrt{2m_2^{-1} \log(2|\mathcal{S}||\mathcal{A}|\delta'^{-1})} \leq (1 - \gamma)u/8. \end{aligned}$$

Finally by shifting the estimate to have one-side error, we obtain the one-side error $(1 - \gamma)u/4$ in the statement of this lemma.

Proof of Lemma 4.4. For i = 0, $Q^{(0)} = r + \gamma w$. By Lemma 4.1, with probability at least $1 - \delta$,

$$|\widetilde{\boldsymbol{w}} - \boldsymbol{P} \boldsymbol{v}^{(0)}| \leq \sqrt{2lpha_1 \boldsymbol{\sigma}_{\boldsymbol{v}^{(0)}}} + rac{2}{3} \cdot lpha_1 \cdot \| \boldsymbol{v}^{(0)} \|_{\infty} \mathbf{1},$$

and

$$\widehat{\boldsymbol{\sigma}} - \boldsymbol{\sigma}_{\boldsymbol{v}^{(0)}} \Big| \le 4 \| \boldsymbol{v}^{(0)} \|_{\infty}^2 \cdot \sqrt{2\alpha_1} \mathbf{1}, \tag{E.3}$$

which we condition on. We have

$$|\widetilde{\boldsymbol{w}} - \boldsymbol{P}\boldsymbol{v}^{(0)}| \leq \sqrt{2\alpha_1}\widehat{\boldsymbol{\sigma}} + (4\alpha_1^{3/4} \|\boldsymbol{v}^{(0)}\|_{\infty} + \frac{2}{3} \cdot \alpha_1 \cdot \|\boldsymbol{v}^{(0)}\|_{\infty})\mathbf{1}.$$

Thus

$$\boldsymbol{w} = \widetilde{\boldsymbol{w}} - \sqrt{2\alpha_1 \widehat{\boldsymbol{\sigma}}} - 4\alpha_1^{3/4} \|\boldsymbol{v}^{(0)}\|_{\infty} \mathbf{1} - \frac{2}{3} \cdot \alpha_1 \cdot \|\boldsymbol{v}^{(0)}\|_{\infty} \mathbf{1} \le \boldsymbol{P} \boldsymbol{v}^{(0)}, \quad (E.4)$$

and

$$\boldsymbol{w} \geq \boldsymbol{P}\boldsymbol{v}^{(0)} - 2\sqrt{2\alpha_1}\widehat{\boldsymbol{\sigma}} - (8\alpha_1^{3/4} \|\boldsymbol{v}^{(0)}\|_{\infty} + \frac{4}{3} \cdot \alpha_1 \cdot \|\boldsymbol{v}^{(0)}\|_{\infty}) \mathbf{1}.$$

By (E.3) and Lemma 4.2, we have

$$\sqrt{\hat{\sigma}} \le \sqrt{\sigma_{v^{(0)}}} + 2\|v^{(0)}\|_{\infty} (2\alpha)^{1/4} \mathbf{1} \le \sqrt{\sigma_{v^*}} + u\mathbf{1} + 2\|v^{(0)}\|_{\infty} (2\alpha)^{1/4} \mathbf{1}$$

we have

$$\boldsymbol{w} \ge \boldsymbol{P}\boldsymbol{v}^{(0)} - 2\sqrt{2\alpha_1}\boldsymbol{\sigma}_{\boldsymbol{v}^*} - 2\sqrt{2\alpha_1}u\boldsymbol{1} - 16\alpha_1^{3/4} \|\boldsymbol{v}^{(0)}\|_{\infty}\boldsymbol{1} - \frac{4}{3} \cdot \alpha_1 \cdot \|\boldsymbol{v}^{(0)}\|_{\infty}\boldsymbol{1}$$
(E.5)

For the rest of the proof, we condition on the event that (E.4) and (E.5) hold, which happens with probability at least $1 - \delta$. Denote $\boldsymbol{v}^{(-1)} = \boldsymbol{0}$. Thus we have $\boldsymbol{v}^{(-1)} \leq \boldsymbol{v}^{(0)} \leq \mathcal{T}_{\pi^{(0)}}(\boldsymbol{v}^{(0)})$. Next we prove the lemma by induction on *i*. Assume for some $i \geq 1$, with probability at least $1 - (i - 1)\delta'$ the following holds,

$$\forall 0 \leq k \leq i-1: \quad \boldsymbol{v}^{(k-1)} \leq \boldsymbol{v}^{(k)} \leq \mathcal{T}_{\pi^{(k)}}(\boldsymbol{v}^{(k)}),$$

which we condition on. Next we show that the lemma statement holds for k = i. By definition of $v^{(i)}$ (Line 27 and 28),

$$oldsymbol{v}^{(i-1)} \leq oldsymbol{v}^{(i)}$$
 and $oldsymbol{v}oldsymbol{(Q}^{(i-1)}ildsymbol) \leq oldsymbol{v}^{(i)}$

Furthermore, since $v^{(0)} \le v^{(1)} \le \dots \le v^{(i-1)} \le \mathcal{T}_{\pi^{i-1}} v^{(i-1)} \le \mathcal{T} v^{(i-1)} \le \mathcal{T}^{\infty} v^{(i-1)} = v^*$, we have

$$v^{(i)} - v^{(0)} \le v^* - v^{(0)} \le u1$$

By Lemma 4.3, we have, with probability at least $1 - \delta'$

$$\boldsymbol{P}[\boldsymbol{v}^{(i)} - \boldsymbol{v}^{(0)}] - \frac{(1 - \gamma)u}{8} \cdot \mathbf{1} \le \boldsymbol{g}^{(i)} \le \boldsymbol{P}[\boldsymbol{v}^{(i)} - \boldsymbol{v}^{(0)}], \quad (E.6)$$

which we condition on for the rest of the proof. Thus we have

$$Q^{(i)} = r + \gamma(w + g^{(i)}) \le r + \gamma(Pv^{(0)} + Pv^{(i)} - Pv^{(0)}) = r + \gamma Pv^{(i)}.$$

To show $\boldsymbol{v}^{(i)} \leq \mathcal{T}_{\pi^{(i)}} \boldsymbol{v}^{(i)}$, we notice that if for some $s, \pi^{(i)}(s) \neq \pi^{(i-1)}(s)$, then

$$\boldsymbol{v}^{(i)}(s) \leq [\mathcal{T}_{\pi^{(i)}} \boldsymbol{v}^{(i-1)}](s) \leq [\mathcal{T}_{\pi^{(i)}} \boldsymbol{v}^{(i)}](s),$$

where the first inequality follows from $\boldsymbol{v}^{(i)}(s) \leq \boldsymbol{r}(s, \pi^{(i)}(s)) + \gamma \boldsymbol{P}_{s,\pi^{(i)}(s)}^{\top} \boldsymbol{v}^{(i-1)} = \mathcal{T}_{\pi^{(i)}} \boldsymbol{v}^{(i-1)}$. On the other hand, if $\pi^{(i)}(s) = \pi^{(i-1)}(s)$, then

$$\boldsymbol{v}^{(i)}(s) = \boldsymbol{v}^{(i-1)}(s) \le (\mathcal{T}_{\pi^{(i-1)}} \boldsymbol{v}^{(i-1)})(s) \le (\mathcal{T}_{\pi^{(i-1)}} \boldsymbol{v}^{(i)})(s) = (\mathcal{T}_{\pi^{(i)}} \boldsymbol{v}^{(i)})(s).$$

This completes the induction step. Lastly, combining (E.5) and (E.6), we have

$$\begin{aligned} \boldsymbol{Q}^{*} - \boldsymbol{Q}^{(i)} &= \boldsymbol{Q}^{*} - \boldsymbol{r} - \gamma(\boldsymbol{w} + \boldsymbol{g}^{(i)}) = \gamma \boldsymbol{P} \boldsymbol{v}(\boldsymbol{Q}^{*}) - \gamma(\boldsymbol{w} + \boldsymbol{g}^{(i)}) \\ &= \gamma \boldsymbol{P} \boldsymbol{v}(\boldsymbol{Q}^{*}) - \gamma \boldsymbol{P}(\boldsymbol{v}^{(i)} - \boldsymbol{v}^{(0)}) - \gamma \boldsymbol{P} \boldsymbol{v}^{(0)} + \boldsymbol{\xi}^{(i)} \\ &= \gamma \boldsymbol{P} \boldsymbol{v}(\boldsymbol{Q}^{*}) - \gamma \boldsymbol{P} \boldsymbol{v}^{(i)} + \boldsymbol{\xi}^{(i)}, \end{aligned}$$

where

$$\boldsymbol{\xi}^{(i)} \leq (1-\gamma)u/8 \cdot \mathbf{1} + 2\sqrt{2\alpha_1 \boldsymbol{\sigma}_{\boldsymbol{v}^*}} + 2\sqrt{2\alpha_1}u \cdot \mathbf{1} + 16\alpha_1^{3/4} \|\boldsymbol{v}^{(0)}\|_{\infty} \cdot \mathbf{1} + (4/3) \cdot \alpha_1 \cdot \|\boldsymbol{v}^{(0)}\|_{\infty} \cdot \mathbf{1},$$

where $\alpha_1 = \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})/m_1 \leq 1$. Mover, since $\boldsymbol{v}(\boldsymbol{Q}^{(i-1)}) \leq \boldsymbol{v}^{(i)}$, we obtain

$$\boldsymbol{Q}^* - \boldsymbol{Q}^{(i)} \leq \gamma \boldsymbol{P} \boldsymbol{v}(\boldsymbol{Q}^*) - \gamma \boldsymbol{P} \boldsymbol{v}(\boldsymbol{Q}^{(i-1)}) + \boldsymbol{\xi}^{(i)} \leq \gamma \boldsymbol{P}^{\pi^*} \boldsymbol{Q}^* - \gamma \boldsymbol{P}^{\pi^*} \boldsymbol{Q}^{(i-1)} + \boldsymbol{\xi}^{(i)},$$

where π^* is an arbitrary optimal policy and we use the fact that $\max_a \mathbf{Q}^*(s, a) = \mathbf{Q}^*(s, \pi^*(s))$. This completes the proof of the lemma. Proof of Proposition 4.5. Recall that we are able to sample a state from each $P_{s,a}$ with time O(1). Let $\beta = (1 - \gamma)^{-1}$, $R = \lceil c_1\beta \ln[\beta u^{-1}] \rceil$, $m_1 = c_2\beta^3 u^{-2} \cdot \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})$ and $m_2 = c_3\beta^2 \cdot \log[2R|\mathcal{S}||\mathcal{A}|\delta^{-1}]$ for some constants c_1, c_2 and c_3 required in Algorithm 1. In the following proof, we set c_1, c_2, c_3 to be sufficiently large but otherwise arbitrary absolute constants (e.g., $c_1 \ge 4, c_2 \ge 8192, c_3 \ge 128$). By Lemma 4.4, with probability at least $1 - 2\delta$ for each $1 \le i \le R$, we have $v^{(i-1)} \le v^{(i)} \le T_{\pi^{(i)}} v^{(i)}$, and $Q^{(i)} \le r + \gamma P v^{(i)}$,

$$\boldsymbol{Q}^* - \boldsymbol{Q}^{(i)} \leq \gamma \boldsymbol{P}^{\pi^*} \left[\boldsymbol{Q}^* - \boldsymbol{Q}^{(i-1)} \right] + \boldsymbol{\xi},$$

where

$$\boldsymbol{\xi} \leq (1 - \gamma) u/C \cdot \boldsymbol{1} + C \sqrt{\alpha_1 \boldsymbol{\sigma}_{\boldsymbol{v}^*}} + C \alpha_1^{3/4} \| \boldsymbol{v}^{(0)} \|_{\infty} \cdot \boldsymbol{1}$$

for $\alpha_1 = \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})/m_1$ and sufficiently large constant C. Solving the recursion, we obtain

$$oldsymbol{Q}^{*} - oldsymbol{Q}^{(R-1)} \leq \gamma^{R-1} oldsymbol{P}^{\pi^{*}} ig[oldsymbol{Q}^{*} - oldsymbol{Q}^{0} ig] + \sum_{i=0}^{R-1} \gamma^{i} (oldsymbol{P}^{\pi^{*}})^{i} oldsymbol{\xi} \ \leq \gamma^{R-1} oldsymbol{P}^{\pi^{*}} ig[oldsymbol{Q}^{*} - oldsymbol{Q}^{0} ig] + (I - \gamma oldsymbol{P}^{\pi^{*}})^{-1} oldsymbol{\xi}.$$

We first apply a naïve bound $\| \boldsymbol{P}^{\pi^*} [\boldsymbol{Q}^* - \boldsymbol{Q}^0] \|_{\infty} \leq (1 - \gamma)^{-1}$. Hence

$$\gamma^{R-1} \boldsymbol{P}^{\pi^*} \left[\boldsymbol{Q}^* - \boldsymbol{Q}^0
ight] \leq rac{u}{4} \cdot \mathbf{1},$$

where $R = \lceil (1-\gamma)^{-1} \ln[4(1-\gamma)^{-1}u^{-1}] \rceil + 1$. The next step is the key to the improvement in our analysis. We further apply the bound in Lemma C.1, given by

$$(I - \gamma P^{\pi^*})^{-1} \sqrt{\sigma_{v^*}} \le \min(2\gamma^{-1}(1-\gamma)^{-1.5}, (1-\gamma)^{-2}) \cdot \mathbf{1} \le 3(1-\gamma)^{-1.5} \cdot \mathbf{1},$$

where the last inequality follows since $\min(2\gamma^{-1}, (1-\gamma)^{-1/2}) \leq 3$. With $\|(\boldsymbol{I} - \gamma \boldsymbol{P}^{\pi^*})^{-1} \boldsymbol{1}\|_{\infty} \leq (1-\gamma)^{-1}$ and $\|\boldsymbol{v}^{(0)}\|_{\infty} \leq (1-\gamma)^{-1}$, we have,

$$(\mathbf{I} - \gamma \mathbf{P}^{\pi^*})^{-1} \boldsymbol{\xi} \leq \left[\frac{u}{8} + C' \sqrt{\frac{2\alpha_1}{\gamma^2 (1 - \gamma)^3}} + C' \frac{\alpha_1^{3/4}}{(1 - \gamma)^2} \right] \cdot \mathbf{1}$$
$$\leq \left[\frac{u}{8} + \frac{u}{16} + \left(\frac{(1 - \gamma)^3 u^2}{C'' \cdot (1 - \gamma)^{8/3}} \right)^{3/4} \right] \cdot \mathbf{1}$$
$$\leq \frac{u}{4} \cdot \mathbf{1},$$

for some sufficiently large C' and C'', which depend on c_1, c_2 and c_3 . Since $v(Q^{(R-1)}) \leq v^{(R)}$, we have

$$v^* - v^{(R)} \le v^* - v(Q^{(R-1)}) \le \gamma^{R-1} P^{\pi^*} [Q^* - Q^0] + (I - \gamma P^{\pi^*})^{-1} \xi \le \frac{u}{2} \cdot 1.$$

This completes the proof of the correctness. It remains to bound the time complexity. The initialization stage costs $O(m_1)$ time per (s, a). Each iteration costs $O(m_2)$ time per (s, a). We thus have the total time complexity as

$$O(m_1 + Rm_2)|\mathcal{S}|||\mathcal{A}|| = O\left[\frac{|\mathcal{S}||\mathcal{A}|}{(1-\gamma)^3} \cdot \log\frac{|\mathcal{S}||\mathcal{A}|}{\delta \cdot (1-\gamma) \cdot u} \cdot \left(\frac{1}{u^2} + \log\frac{1}{(1-\gamma) \cdot u}\right)\right].$$

Since $\log[(1-\gamma)^{-1}u^{-1}] = O(\log[(1-\gamma)^{-1}]u^{-2})$, we conclude the proof.

E.2 Missing Analysis of Halving Errors

We refer in this section to Algorithm 1 as a subroutine HALFERR, which given an input MDP \mathcal{M} with a sampling oracle, an input value function $v^{(i)}$ and an input policy $\pi^{(i)}$, outputs an value function $v^{(i+1)}$ and a policy $\pi^{(i+1)}$ such that, with high probability (over the new samples of the sampling oracle),

$$\| oldsymbol{Q}^{(i+1)} - oldsymbol{Q}^* \|_{\infty} \le \| oldsymbol{Q}^{(i)} - oldsymbol{Q}^* \|_{\infty} / 2 \quad ext{and} \quad \| oldsymbol{v}^{\pi^{(i+1)}} - oldsymbol{v}^* \|_{\infty} \le \| oldsymbol{v}^{\pi^{(i)}} - oldsymbol{v}^* \|_{\infty} / 2.$$

After $\log[\epsilon^{-1}(1-\gamma)^{-1}]$ calls of the subroutine HALFERR, the final output policy and value functions are ϵ -close to the optimal ones with high probability.

We summarize our meta algorithm in Algorithm 2. Note that in the algorithm, each call of HALFERR will draw new samples from the sampling oracle. These new samples guarantee the independence of successive improvements and also save space of the algorithm. For instance, the algorithm HAL-FERR only needs to use O(|S||A|) words of memory instead of storing all the samples. The guarantee of the algorithm is summarized in Proposition E.3.

Proposition E.3. Let $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \boldsymbol{r}, \boldsymbol{P}, \gamma)$ with a sampling oracle. Suppose HALFERR is an algorithm that takes an input $\boldsymbol{v}^{(i)}$ and an input policy $\pi^{(i)}$ and a number $u \in [0, (1-\gamma)^{-1}]$ satisfying $\boldsymbol{v}^* - u\mathbf{1} \leq \boldsymbol{v}^{(i)} \leq \boldsymbol{v}^{\pi^{(i)}}$, halts in time τ and outputs a $\boldsymbol{v}^{(i+1)}$ and a policy $\pi^{(i+1)}$ satisfying,

$$\boldsymbol{v}^* - rac{u}{2} \cdot \mathbf{1} \leq \boldsymbol{v}^{(i+1)} \leq \boldsymbol{v}^{\pi^{(i+1)}} \leq \boldsymbol{v}^*.$$

with probability at least $1 - (1 - \gamma) \cdot \epsilon \cdot \delta$ (over the randomness of the new samples given by the sampling oracle), then the meta algorithm described in Algorithm 2, given input \mathcal{M} and the sampling oracle, halts in $\tau \cdot \log(\epsilon^{-1} \cdot (1 - \gamma)^{-1})$ and outputs an policy $\pi^{(R)}$ such that

$$oldsymbol{v}^* - \epsilon \cdot oldsymbol{1} \leq oldsymbol{v}^{(R)} \leq oldsymbol{v}^{\pi^{(R)}} \leq oldsymbol{v}^*$$
 ,

with probability at least $1 - \delta$ (over the randomness of all samples drawn from the sampling oracle). Moreover, if HALFERR uses space s, then the meta algorithm uses space s + O(|S||A|). If each call of HALFERR takes m samples from the oracle, then the overall samples taken by Algorithm 2 is $m \cdot \log(\epsilon^{-1} \cdot (1 - \gamma)^{-1})$.

The proof of this proposition is a straightforward application of conditional probability.

Proof of Proposition E.3. The proof follows from a straightforward induction. For simplicity, denote $\beta = (1 - \gamma)^{-1}$. In the meta-algorithm, the initialization is $v^{(0)} = 0$ and $\pi^{(0)}$ is an arbitrary policy. Thus $v^* - \beta \cdot \mathbf{1} \leq v^{(0)} \leq v^{\pi^{(0)}}$. By running the meta-algorithm, we obtain a sequence of value functions and policies: $\{v^{(i)}\}_{i=0}^{R}$ and $\{\pi^{(i)}\}_{i=0}^{R}$. Since each call of the HALFERR uses new samples from the oracle, the sequence of value functions and policies satisfies strong Markov property (given $(v^{(i)}, \pi^{(i)}), (v^{(i+1)}, \pi^{(i+1)})$ is independent with $\{(v^{(j)}, \pi^{(j)})\}_{j=0}^{i-1}$). Thus

$$\Pr\left[\boldsymbol{v}^* - 2^{-R}\beta \cdot \mathbf{1} \leq \boldsymbol{v}^{(R)} \leq \boldsymbol{v}^{\pi^{(R)}}\right]$$

$$\geq \prod_{i=1}^{R} \Pr\left[\boldsymbol{v}^* - 2^{-i}\beta \cdot \mathbf{1} \leq \boldsymbol{v}^{(i)} \leq \boldsymbol{v}^{\pi^{(i)}} \middle| \boldsymbol{v}^* - 2^{-i+1}\beta \cdot \mathbf{1} \leq \boldsymbol{v}^{(i-1)} \leq \boldsymbol{v}^{\pi^{(i-1)}}\right]$$

$$\geq 1 - \delta.$$

Since $2^{-R}(1-\gamma)^{-1} \leq \epsilon$, we conclude the proof.

Proof of Theorem 4.6. Our algorithm is simply plugging in Algorithm 1 as the HALFERR subroutine in Algorithm 2. The correctness is guaranteed by Proposition E.3 and Proposition 4.5. The running time guarantee follows from a straightforward calculation.

F Extension to Finite Horizon

In this section we show how to apply similar techniques to achieve improved sample complexities for solving finite Horizon MDPs given a generative model and we prove that the sample complexity we achieve is optimal up to logarithmic factors.

The finite horizon problem is to compute an optimal non-stationary policy over a fixed time horizon H, i.e. a policy of the form $\pi(s,h)$ for $s \in S$ and $h \in \{0, \ldots H\}$), where the reward is the expected cumulative (un-discounted) reward for following this policy. In classic value iteration, this is typically done using a backward recursion from time $H, H - 1, \ldots 0$. We show how to use the ideas in this paper to solve for an ϵ -approximate policy. As we have shown in the discounted case, it is suffice to show an algorithm that decrease the error of the value at each stage by half. Our algorithm is presented in Algorithm 3.

To analyze the algorithm, we first provide an analogous lemma of Lemma 4.1,

Lemma F.1 (Empirical Estimation Error). Let \tilde{w}_h and $\hat{\sigma}_h$ be computed in Line 10 of Algorithm 3. Recall that \tilde{w}_h and $\hat{\sigma}_h$ are empirical estimates of Pv_h and $\sigma_{v_h} = Pv_h^2 - (Pv_h)^2$ using m_1 samples per (s, a) pair. Then with probability at least $1 - \delta$, for $L \stackrel{\text{def}}{=} \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})$ and every $h = 1, 2, \ldots, H$, we have

$$\left|\tilde{\boldsymbol{w}}_{h} - \boldsymbol{P}^{\top} \boldsymbol{v}_{h}^{(0)}\right| \leq \sqrt{2m_{1}^{-1} \boldsymbol{\sigma}_{\boldsymbol{v}_{h}^{(0)}} \cdot L} + 2(3m_{1})^{-1} \|\boldsymbol{v}_{h}^{(0)}\|_{\infty} L$$
(F.1)

and

$$\forall (s,a) \in \mathcal{S} \times \mathcal{A} : \quad \left| \widehat{\boldsymbol{\sigma}}_h(s,a) - \boldsymbol{\sigma}_{\boldsymbol{v}_h^{(0)}}(s,a) \right| \le 4 \|\boldsymbol{v}_h^{(0)}\|_{\infty}^2 \cdot \sqrt{2m_1^{-1}L}.$$
(F.2)

Proof. The proof of this lemma is identical to that of Lemma 4.1.

An analogous lemma to Lemma 4.3 is also presented here.

Lemma F.2. Let $\boldsymbol{g}_h^{(i)}$ be the estimate of $\boldsymbol{P}[\boldsymbol{v}_h^{(i)} - \boldsymbol{v}_h^{(0)}]$ defined in Line 27 of Algorithm 3. Then conditioning on the event that $\|\boldsymbol{v}_h^{(i)} - \boldsymbol{v}_h^{(0)}\|_{\infty} \leq 2\epsilon$, with probability at least $1 - \delta/H$,

$$\boldsymbol{P}[\boldsymbol{v}_{h}^{(i)} - \boldsymbol{v}_{h}^{(0)}] - \frac{\epsilon}{4H} \cdot \boldsymbol{1} \le \boldsymbol{g}_{h}^{(i)} \le \boldsymbol{P}[\boldsymbol{v}_{h}^{(i)} - \boldsymbol{v}_{h}^{(0)}]$$

provided appropriately chosen constants in Algorithm 3.

Proof. The proof of this lemma is identical to that of Lemma 4.3 except that $(1 - \gamma)^{-1}$ is replaced with *H*.

Similarly, we can show the following improvement lemma.

Lemma F.3. Let Q_h be the estimated Q-function of v_{h+1} in Line 30 of Algorithm 3. Let $Q_h^* = r + P_h v_{h+1}^*$ be the optimal Q-function of the DMDP. Let $\pi(\cdot, h)$ and v_h be estimated in iteration h, as defined in Line 24 and 25. Let π^* be an optimal policy for the DMDP. For a policy π , let $P_h^{\pi}Q \in \mathbb{R}^{S \times A}$ be defined as $(P_h^{\pi}Q)(s,a) = \sum_{s' \in S} P_{s,a}(s')Q(s', \pi(s', h))$. Suppose for all $h \in [H-1]$, $v_h^{(0)} \leq \mathcal{T}_{\pi^{(0)}(\cdot,h)}v_{h+1}^{(0)}$. Let $v_{H+1} \stackrel{\text{def}}{=} 0$ and $Q_{H+1} \stackrel{\text{def}}{=} 0$. Then, with probability at least $1-2\delta$, for all $1 \leq h \leq H$, $v_h^{(0)} \leq v_h \leq \mathcal{T}_{\pi(\cdot,h)}v_{h+1} \leq v_h^*$, $Q_h \leq r + P_hv_{h+1}$ and

$$\boldsymbol{Q}_{h}^{*} - \boldsymbol{Q}_{h} \leq \boldsymbol{P}_{h}^{\pi^{*}} \left[\boldsymbol{Q}_{h+1}^{*} - \boldsymbol{Q}_{h+1} \right] + \boldsymbol{\xi}_{h},$$

where the error vector $\boldsymbol{\xi}_h$ satisfies

$$\mathbf{0} \leq \boldsymbol{\xi}_{h} \leq 8H^{-1}u \cdot \mathbf{1} + 2\sqrt{2\alpha_{1}\boldsymbol{\sigma}_{\boldsymbol{v}_{h+1}^{*}}} + 2\sqrt{2\alpha_{1}}u \cdot \mathbf{1} + 16\alpha_{1}^{3/4} \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \cdot \mathbf{1} + (4/3) \cdot \alpha_{1} \cdot \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \cdot \mathbf{1},$$

and $\alpha_{1} = \log(8|\mathcal{S}||\mathcal{A}|H\delta^{-1})/m_{1}.$

Proof of Lemma F.3. By Lemma 4.1, for any h = 1, 2, ..., H, with probability at least $1 - \delta/H$,

$$|\widetilde{\boldsymbol{w}}_{h} - \boldsymbol{P}\boldsymbol{v}_{h+1}| \leq \sqrt{2\alpha_1 \boldsymbol{\sigma}_{\boldsymbol{v}_{h+1}^{(0)}}} + \frac{2}{3} \cdot \alpha_1 \cdot \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \cdot \mathbf{1},$$

and

$$\left|\widehat{\boldsymbol{\sigma}}_{h+1} - \boldsymbol{\sigma}_{\boldsymbol{v}_{h+1}^{(0)}}\right| \le 4 \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty}^2 \cdot \sqrt{2\alpha_1} \cdot \mathbf{1},\tag{F.3}$$

which we condition on. We have

$$|\widetilde{\boldsymbol{w}}_{h} - \boldsymbol{P}\boldsymbol{v}_{h+1}^{(0)}| \le \sqrt{2\alpha_{1}\widehat{\boldsymbol{\sigma}}_{h+1}} + (4\alpha_{1}^{3/4} \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} + \frac{2}{3} \cdot \alpha_{1} \cdot \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty})\mathbf{1}.$$

Thus

$$\boldsymbol{w}_{h} = \widetilde{\boldsymbol{w}}_{h} - \sqrt{2\alpha_{1}\widehat{\boldsymbol{\sigma}}_{h+1}} - 4\alpha_{1}^{3/4} \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \mathbf{1} - \frac{2}{3} \cdot \alpha_{1} \cdot \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \mathbf{1} \le \boldsymbol{P}\boldsymbol{v}_{h+1}^{(0)}, \quad (F.4)$$

and

$$m{w}_h \ge m{P}m{v}_{h+1}^{(0)} - 2\sqrt{2lpha_1}\widehat{m{\sigma}}_{h+1} - (8lpha_1^{3/4} \|m{v}_{h+1}^{(0)}\|_\infty + rac{4}{3} \cdot lpha_1 \cdot \|m{v}_{h+1}^{(0)}\|_\infty) \mathbf{1}.$$

By (E.3) and Lemma 4.2, we have

$$\sqrt{\hat{\sigma}_{h+1}} \le \sqrt{\sigma_{\boldsymbol{v}_{h+1}^{(0)}}} + 2\|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} (2\alpha)^{1/4} \mathbf{1} \le \sqrt{\sigma_{\boldsymbol{v}_{h+1}^*}} + \epsilon \mathbf{1} + 2\|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} (2\alpha)^{1/4} \mathbf{1}.$$

we have

$$\boldsymbol{w}_{h} \geq \boldsymbol{P}\boldsymbol{v}_{h+1}^{(0)} - 2\sqrt{2\alpha_{1}\boldsymbol{\sigma}_{\boldsymbol{v}_{h+1}^{*}}} - 2\sqrt{2\alpha_{1}}\epsilon \boldsymbol{1} - 16\alpha_{1}^{3/4} \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \boldsymbol{1} - \frac{4}{3} \cdot \alpha_{1} \cdot \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \boldsymbol{1}$$
(F.5)

For the rest of the proof, we condition on the event that (F.4) and (F.5) hold for all h = 1, 2, ..., H, which happens with probability at least $1 - \delta$. Denote $v_{H+1}^* = v_{H+1} = v_{H+1}^{(0)} = 0$. Thus we have $v_{H+1}^{(0)} \leq v_{H+1} \leq v_{H+1}^*$. Next we prove the lemma by induction on h. Assume for some h, with probability at least $1 - (h - 1)\delta/H$ the following holds, for all h' = h + 1, h + 2, ..., H,

$$oldsymbol{v}_{h'}^{(0)} \leq oldsymbol{v}_{h'} \leq oldsymbol{v}_{h'}^*$$

which we condition on. Next we show that the lemma statement holds for h as well. By definition of v_h (Line 27 and 28),

$$oldsymbol{v}_h^{(0)} \leq oldsymbol{v}_h$$
 and $oldsymbol{v}(oldsymbol{Q}_h) \leq oldsymbol{v}_h.$

Furthermore, since $\boldsymbol{v}_{h+1}^{(0)} \leq \boldsymbol{v}_{h+1}^* \leq \boldsymbol{v}_{h+1}^{(0)} + u \mathbf{1}$ we have

$$v_{h+1}^* - v_{h+1} \le v_{h+1}^* - v_{h+1}^{(0)} \le u\mathbf{1}.$$

By Lemma 4.3, we have, with probability at least $1 - \delta'$

$$\boldsymbol{P}[\boldsymbol{v}_{h+1} - \boldsymbol{v}_{h+1}^{(0)}] - \frac{u}{8H} \cdot \mathbf{1} \le \boldsymbol{g}_h \le \boldsymbol{P}[\boldsymbol{v}_{h+1} - \boldsymbol{v}_{h+1}^{(0)}],$$
(F.6)

which we condition on for the rest of the proof. Thus we have

$$m{Q}_h = m{r} + (m{w}_h + m{g}_h) \le m{r} + m{P}m{v}_{h+1}^{(0)} + m{P}m{v}_{h+1} - m{P}m{v}_{h+1}^{(0)} = m{r} + m{P}m{v}_{h+1} \le m{Q}_h^*.$$

To show $\boldsymbol{v}_h \leq \mathcal{T}_{\pi(\cdot,h)} \boldsymbol{v}_{h+1}$, we notice that if for some $s, \pi(s,h) \neq \pi^{(0)}(s,h)$, then,

$$\boldsymbol{v}_h(s) \leq \boldsymbol{r}(s, \pi(s, h)) + \boldsymbol{P}_{s, \pi(s, h)}^{\top} \boldsymbol{v}_{h+1} = \mathcal{T}_{\pi(\cdot, h)} \boldsymbol{v}_{h+1}.$$

On the other hand, if $\pi(s,h) = \pi^{(0)}(s,h)$, then

$$\forall s \in \mathcal{S}: \quad \boldsymbol{v}_{h}(s) = \boldsymbol{v}_{h}^{(0)}(s) \leq (\mathcal{T}_{\pi^{(0)}(\cdot,h)}\boldsymbol{v}_{h+1}^{(0)})(s) \leq (\mathcal{T}_{\pi^{(0)}(\cdot,h)}\boldsymbol{v}_{h+1})(s) = (\mathcal{T}_{\pi(\cdot,h)}\boldsymbol{v}_{h+1})(s).$$

This completes the induction step. Lastly, combining (F.5) and (F.6), we have

$$egin{aligned} m{Q}_h^* &- m{Q}_h = m{Q}_h^* - m{r} - (m{w}_h + m{g}_h) = m{P}m{v}(m{Q}_{h+1}^*) - (m{w}_h + m{g}_h) \ &= m{P}m{v}(m{Q}_{h+1}^*) - m{P}(m{v}_{h+1} - m{v}_{h+1}^{(0)}) - m{P}m{v}_{h+1} + m{\xi}_h \ &= m{P}m{v}(m{Q}_{h+1}^*) - m{P}m{v}_{h+1} + m{\xi}_h, \end{aligned}$$

where

$$\boldsymbol{\xi}_{h} \leq H^{-1}u/8 \cdot \mathbf{1} + 2\sqrt{2\alpha_{1}\boldsymbol{\sigma}_{\boldsymbol{v}_{h+1}^{*}}} + 2\sqrt{2\alpha_{1}}u \cdot \mathbf{1} + 16\alpha_{1}^{3/4} \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \cdot \mathbf{1} + (4/3) \cdot \alpha_{1} \cdot \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \cdot \mathbf{1},$$

where $\alpha_1 = \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})/m_1$. Mover, since $v(Q_{h+1}) \leq v_{h+1}$, we obtain

$$m{Q}_h^* - m{Q}_h \leq m{P}m{v}(m{Q}_{h+1}^*) - m{P}m{v}(m{Q}_{h+1}) + m{\xi}_h \leq m{P}_h^{\pi^*}m{Q}_{h+1}^* - m{P}_h^{\pi^*}m{Q}_{h+1} + m{\xi}_h$$

where π^* is an arbitrary optimal policy and we use the fact that $\max_a \mathbf{Q}_h^*(s, a) = \mathbf{Q}_h^*(s, \pi^*(s, h))$. This completes the proof of the lemma.

Furthermore, we show an analogous lemma of Lemma C.1.

Lemma F.4 (Upper Bound on Variance). For any π , we have

$$\left\|\sum_{h'=h}^{H-1} \left(\prod_{i=h+1}^{h'} \boldsymbol{P}_i^{\pi}\right) \sqrt{\boldsymbol{\sigma}_{\boldsymbol{v}_{h'+1}}}\right\|_{\infty}^2 \leq H^{3/2}.$$

Proof. First, by Cauchy-Swartz inequality, we have

$$\sum_{h'=h}^{H-1} \bigg(\prod_{i=h+1}^{h'} \boldsymbol{P}_i^{\pi}\bigg) \sqrt{\boldsymbol{\sigma}_{\boldsymbol{v}_{h'+1}}^{\pi}} \leq \sqrt{H \sum_{h'=h}^{H-1} \bigg(\prod_{i=h+1}^{h'} \boldsymbol{P}_i^{\pi}\bigg) \boldsymbol{\sigma}_{\boldsymbol{v}_{h'+1}}^{\pi}}.$$

Next, by a similar argument of the proof of Lemma C.2, we can show that

$$\left[\sum_{h'=h}^{H-1} \left(\prod_{i=h+1}^{h'} \boldsymbol{P}_{i}^{\pi}\right) \boldsymbol{\sigma}_{\boldsymbol{v}_{h'+1}}^{\pi}\right](s) = \operatorname{Var}\left[\sum_{t=h}^{H} r(s^{t}, \pi(s^{t}, t)) \middle| s^{h} = s\right] \le H^{2}.$$

etes the proof.

This completes the proof.

We are now ready to present the guarantee of the algorithm 3.

Proposition F.5. On an input value vectors $\boldsymbol{v}_1^{(0)}, \boldsymbol{v}_2^{(0)}, \dots, \boldsymbol{v}_H^{(0)}$, policy $\pi^{(0)}$, and parameters $u \in (0, \beta], \delta \in (0, 1)$ such that $\boldsymbol{v}_h^{(0)} \leq \mathcal{T}_{\pi^{(0)}(\cdot,h)} \boldsymbol{v}_{h+1}^{(0)}$ for all $h \in [H-1]$, and $\boldsymbol{v}_h^{(0)} \leq \boldsymbol{v}_h^* \leq \boldsymbol{v}_h^{(0)} + u\mathbf{1}$, Algorithm 3 halts in time $O[u^{-2} \cdot H^4|\mathcal{S}||\mathcal{A}| \cdot \log(|\mathcal{S}||\mathcal{A}\delta^{-1}Hu^{-1})]$ and outputs $\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_H$ and $\pi: \mathcal{S} \times [H] \to \mathcal{A}$ such that

$$\forall h \in [H]: \quad \boldsymbol{v}_h \leq \mathcal{T}_{\pi(\cdot,h)}(\boldsymbol{v}_{h+1}) \quad \text{and} \quad \boldsymbol{0} \leq \boldsymbol{v}_h^* - \boldsymbol{v}_h \leq (u/2) \cdot \boldsymbol{1}$$

with probability at least $1-\delta$, provided appropriately chosen constants, c_1 , c_2 and c_3 , in Algorithm 3. Moreover, the algorithm uses $O[u^{-2} \cdot H^3 |S| |A| \cdot \log(|S| |A\delta^{-1} Hu^{-1})]$ samples from the sampling oracle.

Proof of Proposition F.5. Recall that we are able to sample a state from each $P_{s,a}$ with time O(1). Let $R = \lceil c_1 H \ln[Hu^{-1}] \rceil$, $m_1 = c_2 H^3 u^{-2} \cdot \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})$ and $m_2 = c_3 H^2 \cdot \log[2R|\mathcal{S}||\mathcal{A}|\delta^{-1}]$ for some constants c_1, c_2 and c_3 required in Algorithm 1. In the following proof, we set $c_1 = 4, c_2 = 8192, c_3 = 128$. By Lemma 4.4, with probability at least $1 - 2\delta$ for each $1 \leq h \leq H$, we have $v_h^{(0)} \leq v_h \leq \mathcal{T}_{\pi(\cdot,h)}v_h$, and $Q_h \leq r + Pv_{h+1}$,

$$oldsymbol{Q}_h^* - oldsymbol{Q}_h \leq oldsymbol{P}_h^{\pi^*}ig[oldsymbol{Q}_{h+1}^* - oldsymbol{Q}_{h+1}ig] + oldsymbol{\xi}_h,$$

where

$$\boldsymbol{\xi}_{h} \leq H^{-1}u/8 \cdot \mathbf{1} + 2\sqrt{2\alpha_{1}\boldsymbol{\sigma}_{\boldsymbol{v}_{h+1}^{*}}} + 2\sqrt{2\alpha_{1}}u \cdot \mathbf{1} + 16\alpha_{1}^{3/4} \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \cdot \mathbf{1} + (4/3) \cdot \alpha_{1} \cdot \|\boldsymbol{v}_{h+1}^{(0)}\|_{\infty} \cdot \mathbf{1},$$

and $\alpha_1 = \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})/m_1$. Notice that $\boldsymbol{v}_H^{(0)} = \boldsymbol{v}_H^* = \boldsymbol{v}(\boldsymbol{r})$, thus the $\boldsymbol{v}_H - \boldsymbol{v}_H^* = \boldsymbol{0}$. Solving the recursion, we obtain

$$oldsymbol{Q}_h^* - oldsymbol{Q}_h \leq \sum_{h'=h}^{H-1} igg(\prod_{i=h+1}^{h'} oldsymbol{P}_i^{\pi^*}igg) oldsymbol{\xi}_{h'},$$

The next step is the key to the improvement in our analysis. We further apply the bound in Lemma C.1, given by

$$\sum_{h'=h}^{H-1} \left(\prod_{i=h+1}^{h'} \boldsymbol{P}_i^{\pi^*}\right) \sqrt{\boldsymbol{\sigma}_{\boldsymbol{v}_{h'+1}^*}} \le H^{3/2} \cdot \mathbf{1}.$$

With $\|\sum_{h'=h}^{H-1} \prod_{i=h+1}^{h'} \boldsymbol{P}_i^{\pi^*} \mathbf{1}\|_{\infty} \leq H-h+1$ and $\|\boldsymbol{v}_h^{(0)}\|_{\infty} \leq H$, we have,

$$\begin{split} \sum_{h'=h}^{H-1} \left(\prod_{i=h+1}^{h'} \boldsymbol{P}_{i}^{\pi^{*}}\right) \boldsymbol{\xi}_{h'} &\leq \left[\frac{u}{8} + 4\sqrt{2\alpha_{1}H^{3}} + 2H\sqrt{2\alpha_{1}}u + 16H^{2}\alpha_{1}^{3/4} + \frac{4\alpha_{1}H^{2}}{3}\right] \mathbf{1} \\ &\leq \left[\frac{u}{8} + \frac{u}{16} + \frac{\sqrt{H^{-1}u}}{32} + 16\left(\frac{H^{-3}u^{2}}{32 \cdot 256 \cdot (H)^{-8/3}}\right)^{3/4} + \frac{4H^{-1}u^{2}}{24 \cdot 256}\right] \cdot \mathbf{1} \\ &\leq \frac{u}{4} \cdot \mathbf{1}, \end{split}$$

provided

$$\alpha_1 = \frac{\log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})}{m_1} = c_2^{-1}H^3u^{-2} \le \frac{H^{-3}u^2}{32 \cdot 256}$$

Since $\boldsymbol{v}(\boldsymbol{Q}_h) \leq \boldsymbol{v}_h$, we have

$$oldsymbol{v}_h^* - oldsymbol{v}_h \leq oldsymbol{v}^* - oldsymbol{v}(oldsymbol{Q}_h) \leq \sum_{h'=h}^{H-1} igg(\prod_{i=h+1}^{h'} oldsymbol{P}_i^{\pi^*}igg) oldsymbol{\xi}_{h'} \leq rac{u}{2} \cdot oldsymbol{1}.$$

This completes the proof of the correctness. It remains to bound the time complexity. The initialization stage costs $O(m_1)$ time per (s, a) per stage h. Each iteration costs $O(m_2)$ time per (s, a). We thus have the total time complexity as

$$O(Hm_1 + Hm_2)|\mathcal{S}||\mathcal{A}| = O\left[H^4 \cdot |\mathcal{S}||\mathcal{A}| \cdot \log \frac{H|\mathcal{S}||\mathcal{A}|}{\delta \cdot u} \cdot \frac{1}{u^2}\right]$$

The total number of samples used is

$$O(m_1 + Hm_2)|\mathcal{S}||\mathcal{A}| = O\left[H^3 \cdot |\mathcal{S}||\mathcal{A}| \cdot \log \frac{H|\mathcal{S}||\mathcal{A}|}{\delta \cdot u} \cdot \frac{1}{u^2}\right].$$

This completes the proof.

We can then use our meta-algorithm and obtain the following theorem.

Theorem F.6. Let $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r}, H)$ be a H-MDP with a sampling oracle. Suppose we can sample a state from each probability vector $\mathbf{P}_{s,a}$ within time O(1). Then there exists an algorithm that runs in time

$$O\left[\frac{1}{\epsilon^2} \cdot H^4 |\mathcal{S}| |\mathcal{A}| \cdot \log \frac{H|\mathcal{S}| |\mathcal{A}|}{\delta \cdot \epsilon} \cdot \log \frac{H}{\epsilon}\right]$$

and obtains a policy π such that, with probability at least $1 - \delta$,

$$\forall h \in [H] : \boldsymbol{v}_h^* - \epsilon \mathbf{1} \leq \boldsymbol{v}_h^\pi \leq \boldsymbol{v}_h^*,$$

where v_h^* is the optimal value of \mathcal{M} at stage h. Moreover, the number of samples used by the algorithm is

$$O\left[\frac{1}{\epsilon^2} \cdot H^3 |\mathcal{S}| |\mathcal{A}| \cdot \log \frac{H|\mathcal{S}| |\mathcal{A}|}{\delta \cdot \epsilon} \cdot \log \frac{H}{\epsilon}\right].$$

F.1 Sample Lower Bound On *H*-MDP

In this section we show that the sample complexity obtained by the algorithm in the last section is essentially tight. Our proof idea is simple, we will reduce the *H*-MDP problem to a discounted MDP problem. If there is an algorithm that solves an *H*-MDP to obtain an ϵ -optimal value, it also gives an value function to the discounted MDP. Therefore, the lower bound of solving *H*-MDP inherits from that of the discounted MDP. The formal guarantee is presented in the following theorem.

Theorem F.7. Let S and A be finite sets of states and actions. Let H > 0 be a positive integer and $\epsilon \in (0, 1/2)$ be an error parameter. Let K be an algorithm that, on input an H-MDP $\mathcal{M} \stackrel{\text{def}}{=} (S, \mathcal{A}, P, \mathbf{r})$ with a sampling oracle, outputs a value function \mathbf{v}_1 for the first stage, such that $\|\mathbf{v}_1 - \mathbf{v}_1^*\|_{\infty} \leq \epsilon$ with probability at least 0.9. Then K calls the sampling oracle at least $\Omega(H^{-3}\epsilon^{-2}|S||\mathcal{A}|/\log\epsilon^{-1})$ times on some input P and $\mathbf{r} \in [0, 1]^S$.

Proof. Let $s_0 \in S$ be a state. Denote $S' = S \setminus \{s_0\}$ be a subset of S. Let $\gamma \in (0, 1)$ be such that $(1 - \gamma)^{-1} \log \epsilon^{-1} \leq H$. Suppose we have an DMDP $\mathcal{M}' = (S', \mathcal{A}, P', \gamma, \mathbf{r}')$ with a sampling oracle. Let $\mathbf{v}^{*'}$ be the optimal value function of \mathcal{M}' . Note that $\mathbf{v}^{*'} \in \mathbb{R}^{S'}$. We will show, in the next paragraph, an H-MDP $\mathcal{M} = (S, \mathcal{A}, P, H, \mathbf{r})$ with first stage value \mathbf{v}_1^* , such that $\|\mathbf{v}_1^*\|_{S'} - \mathbf{v}^{*'}\| \leq \epsilon$. Therefore, an ϵ -approximation of \mathbf{v}_1^* gives a 2ϵ -approximation to \mathbf{v}^* . We show that \mathcal{K} can be used to obtain an ϵ -approximate value \mathbf{v}_1 for \mathbf{v}_1^* of \mathcal{M} and thus \mathcal{K} inherits the lower bound for obtaining (2ϵ) -approximated value for γ -DMDPs.

For \mathcal{M} , in each state $s \in S'$, for any action there is a $(1 - \gamma)$ probability transiting to s_0 and γ probability to do the original transitions in \mathcal{M}' ; for s_0 , no matter what action taken, it transits to

Algorithm 3 FiniteHorizonRandomQVI

1: Input: $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbf{r}, \mathbf{P})$ with a sampling oracle, $\mathbf{v}_1^{(0)}, \mathbf{v}_2^{(0)}, \dots, \mathbf{v}_H^{(0)}, \pi^{(0)} : \mathcal{S} \times [H] \rightarrow \mathcal{S}$ $\mathcal{A}, u, \delta \in (0, 1);$ 2: $\langle u \rangle$ is the initial error, $\pi^{(0)}$ is the input policy, and δ is the error probability 3: **Output:** $v_1, v_2, ..., v_H, \pi$ $4 \cdot$ 5: INITIALIZATION: 6: Let $m_1 \leftarrow c_1 H^3 u^{-2} \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})$ for constant c_1 ; 7: Let $m_2 \leftarrow c_2 H^2 \log[2H|\mathcal{S}||\mathcal{A}|\delta^{-1}]$ for constant c_2 ; 8: Let $\alpha_1 \leftarrow m_1^{-1} \log(8|\mathcal{S}||\mathcal{A}|\delta^{-1})$; 9: For each $(s, a) \in \mathcal{S} \times \mathcal{A}$, sample independent samples $s_{s,a}^{(1)}, s_{s,a}^{(2)}, \ldots, s_{s,a}^{(m_1)}$ from $P_{s,a}$; 10: Initialize $\boldsymbol{w}_h = \widetilde{\boldsymbol{w}}_h = \widehat{\boldsymbol{\sigma}}_h = \boldsymbol{Q}_h^{(0)} \leftarrow \boldsymbol{0}_{\mathcal{S} \times \mathcal{A}}$ for all $h \in [H]$, and $i \leftarrow 0$; 11: Denote $\boldsymbol{v}_{H+1} \leftarrow \boldsymbol{0}$ and $\boldsymbol{Q}_{H+1} \leftarrow \boldsymbol{0}$ 12: for each $(s, a) \in \mathcal{S} \times \mathcal{A}$, $h \in [H]$ do \\Compute empirical estimates of $P_{s,a}^{\top} v_h^{(0)}$ and $\sigma_{v_s^{(0)}}(s,a)$ 13: Let $\widetilde{\boldsymbol{w}}_{h}(s, a) \leftarrow \frac{1}{m_{1}} \sum_{j=1}^{m_{1}} \boldsymbol{v}_{h}^{(0)}(s_{s,a}^{(j)})$ Let $\widehat{\boldsymbol{\sigma}}_{h}(s, a) \leftarrow \frac{1}{m_{1}} \sum_{j=1}^{m_{1}} (\boldsymbol{v}_{h}^{(0)})^{2}(s_{s,a}^{(j)}) - \widetilde{\boldsymbol{w}}_{h}^{2}(s, a)$ 14: 15: 16: \\Shift the empirical estimate to have one-sided error 17: $\boldsymbol{w}_h(s,a) \leftarrow \widetilde{\boldsymbol{w}}_h(s,a) - \sqrt{2\alpha_1 \widehat{\boldsymbol{\sigma}}_h(s,a)} - 4\alpha_1^{3/4} \|\boldsymbol{v}_h^{(0)}\|_{\infty} - (2/3)\alpha_1 \|\boldsymbol{v}_h^{(0)}\|_{\infty}$ 18: 19: Let $v_{H+1} \leftarrow 0$ and $Q_{H+1} \leftarrow 0$. 20: 21: **REPEAT:** \\successively improve 22: for h = H, H - 1 to 1 do $\begin{array}{l} \langle \mathsf{Compute} \ \boldsymbol{P}_{s,a}^{\top} \left[\boldsymbol{v}_{h} - \boldsymbol{v}_{h}^{(0)} \right] \text{ with one-sided error} \\ \text{Let } \widetilde{\boldsymbol{v}}_{h} \leftarrow \boldsymbol{v}_{h} \leftarrow \boldsymbol{v}(\boldsymbol{Q}_{h+1}), \widetilde{\pi}(\cdot,h) \leftarrow \pi(\cdot,h) \leftarrow \pi(\boldsymbol{Q}_{h+1}), \boldsymbol{v}_{h} \leftarrow \widetilde{\boldsymbol{v}}_{h}; \\ \text{For each } s \in \mathcal{S}, \text{ if } \widetilde{\boldsymbol{v}}_{h}(s) \leq \boldsymbol{v}_{h}^{(0)}(s), \text{ then } \boldsymbol{v}_{h}(s) \leftarrow \boldsymbol{v}_{h}^{(0)}(s) \text{ and } \pi(s,h) \leftarrow \pi^{(0)}(s,h); \end{array}$ 23: 24: 25: For each $(s, a) \in \mathcal{S} \times \mathcal{A}$, sample independent samples $\widetilde{s}_{s,a}^{(1)}, \widetilde{s}_{s,a}^{(2)}, \ldots, \widetilde{s}_{s,a}^{(m_2)}$ from $P_{s,a}$; 26: Let $\boldsymbol{g}_h(s,a) \leftarrow m_2^{-1} \sum_{i=1}^{m_2} \left[\boldsymbol{v}_h(\widetilde{s}_{s,a}^{(j)}) - \boldsymbol{v}_h^{(0)}(\widetilde{s}_{s,a}^{(j)}) \right] - H^{-1} u/8;$ 27: 28: $\backslash \backslash Improve Q_h$: 29: $Q_h \leftarrow r + w_h + g_h;$ 30: 31: return $v_1, v_2, \ldots, v_H, \pi$.

itself with probability 1. Formally, for each state $s, s' \in \mathcal{S}', a \in \mathcal{A}$, $P(s'|s, a) = \gamma \cdot P'(\cdot|s, a)$ and $P(s_0|s, a) = (1 - \gamma)$; $P(s'|s_0, a) = 0$ and $P(s_0|s_0, a) = 1$. For \boldsymbol{r} , we set $\boldsymbol{r}(s_0, \cdot) = \boldsymbol{0}$ and $\boldsymbol{r}(s, \cdot) = \boldsymbol{r}'(s, \cdot)$ for $s \in \mathcal{S}'$. It remains to show that $\|\boldsymbol{v}_1^*\|_{\mathcal{S}'} - \boldsymbol{v}^*\|_{\infty} \leq \epsilon$. First we note that $\boldsymbol{v}(\boldsymbol{r}) = \boldsymbol{v}_H^* \leq \boldsymbol{v}^*$. Then, by monotonicity of the \mathcal{T} operator, we have, for all $h \in [H - 1]$ and $s \in \mathcal{S}'$,

$$\boldsymbol{v}_h^*|_{\mathcal{S}'}(s) = \max_{\boldsymbol{a}} [\boldsymbol{r}'(s, \boldsymbol{a}) + \gamma \boldsymbol{P}_{s, \boldsymbol{a}}^{'\top} \boldsymbol{v}_{h+1}^*] \leq \boldsymbol{v}^{*'}.$$

In particular, $v_1^*|_{S'} \leq v^{*'}$. Since the optimal policy $\pi^{*'}$ of \mathcal{M}' can be used as a policy for the *H*-MDP as a non-optimal one, we have

$$oldsymbol{v}^* - \epsilon \cdot oldsymbol{1} \leq igg[1 + \gamma oldsymbol{P}_{\pi^{*'}} + \gamma^2 oldsymbol{P}_{\pi^{*'}}^2 + \cdot + \gamma^H \cdot oldsymbol{P}_{\pi^{*'}}^Higg] oldsymbol{r}^{\pi^{*'}} \leq oldsymbol{v}_1^*|_{\mathcal{S}'}.$$

This completes the proof.

The above lower bound with our algorithm also implies a sample lower bound for an ϵ -policy.

Corollary F.8. Let S and A be finite sets of states and actions. Let H > 0 be a positive integer and $\epsilon \in (0, 1/2)$ be an error parameter. Let K be an algorithm that, on input an H-MDP

 $\mathcal{M} := (\mathcal{S}, \mathcal{A}, P, \mathbf{r})$ with a sampling oracle, outputs a policy $\pi : \mathcal{S} \times [H] \to \mathcal{A}$, such that $\forall h : \|\mathbf{v}_h^{\pi} - \mathbf{v}_h^*\|_{\infty} \leq \epsilon$ with probability at least 0.9. Then \mathcal{K} calls the sampling oracle at least $\Omega(H^{-3}\epsilon^{-2}|\mathcal{S}||\mathcal{A}|/\log\epsilon^{-1})$ times on the worst case input P and $\mathbf{r} \in [0, 1]^{\mathcal{S}}$.