
Efficiently Solving MDPs with Stochastic Mirror Descent

Yujia Jin¹ Aaron Sidford¹

Abstract

We present a unified framework based on primal-dual stochastic mirror descent for approximately solving infinite-horizon Markov decision processes (MDPs) given a generative model. When applied to an average-reward MDP with A_{tot} total actions and mixing time bound t_{mix} our method computes an ϵ -optimal policy with an expected $\tilde{O}(t_{\text{mix}}^2 A_{\text{tot}} \epsilon^{-2})$ samples from the state-transition matrix, removing the ergodicity dependence of prior art. When applied to a γ -discounted MDP with A_{tot} total actions our method computes an ϵ -optimal policy with an expected $\tilde{O}((1-\gamma)^{-4} A_{\text{tot}} \epsilon^{-2})$ samples, improving over previous primal-dual methods and matching the state-of-the-art up to a $(1-\gamma)^{-1}$ factor. Both methods are model-free, update state values and policies simultaneously, and run in time linear in the number of samples taken. We achieve these results through a more general stochastic mirror descent framework for solving bilinear saddle-point problems with simplex and box domains and we demonstrate the flexibility of this framework by providing further applications to constrained MDPs.

1 Introduction

Markov decision processes (MDPs) are a fundamental mathematical abstraction for sequential decision making under uncertainty and they serve as a basic modeling tool in reinforcement learning (RL) and stochastic control (Bertsekas & Tsitsiklis, 1995; Puterman, 2014; Sutton & Barto, 2018). Two prominent classes of MDPs are discounted MDPs (DMDPs) and average-reward MDPs (AMDPs). Each have been studied extensively; DMDPs have a number of nice theoretical properties including reward convergence and operator monotonicity (Bertsekas et al., 1995) and AMDPs are

¹Management Science and Engineering (MS&E), Stanford University, Stanford, United States. Correspondence to: Yujia Jin <yujiajin@stanford.edu>.

applicable to optimal control, learning automata, and various real-world reinforcement learning settings (Mahadevan, 1996; Auer & Ortner, 2007; Ouyang et al., 2017).

In this paper we consider the prevalent computational learning problem of finding an approximately optimal policy of an MDP given only restricted access to the model. In particular, we consider the problem of computing an ϵ -optimal policy, i.e. a policy with an additive ϵ error in expected cumulative reward over infinite horizon, under the standard assumption of a generative model (Kearns & Singh, 1999; Kakade et al., 2003), which allows one to sample from state-transitions given the current state-action pair. This problem is well-studied and there are multiple known upper and lower bounds on its sample complexity (Azar et al., 2012; Wang, 2017a; Sidford et al., 2018a; Wainwright, 2019).

In this work, we provide a unified framework based on primal-dual stochastic mirror descent (SMD) for learning an ϵ -optimal policies for both AMDPs and DMDPs with a generative model. We show that this framework achieves sublinear running times for solving dense bilinear saddle-point problems with simplex and box domains, and (as a special case) ℓ_∞ regression (Sherman, 2017; Sidford & Tian, 2018). As far as we are aware, this is the first such sub-linear running time for this problem. We achieve our results by applying this framework to saddle-point representations of AMDPs and DMDPs and proving that approximate equilibria yield approximately optimal policies.

Our MDP algorithms have sample complexity linear in the total number of possible actions, denoted by A_{tot} . For an AMDP with bounded mixing time t_{mix} for all policies, we prove a sample complexity of $\tilde{O}(t_{\text{mix}}^2 A_{\text{tot}} \epsilon^{-2})$ ¹, which removes the ergodicity condition of prior art (Wang, 2017b) (which can in the worst-case be unbounded). For DMDP with discount factor γ , we prove a sample complexity of $\tilde{O}((1-\gamma)^{-4} A_{\text{tot}} \epsilon^{-2})$, improving over the best-known achievable sample complexity by primal-dual methods (Wang, 2017a) and matching the state-of-the-art (Sidford et al., 2018a; Wainwright, 2019) and lower bound (Azar et al., 2012) up to a $(1-\gamma)^{-1}$ factor.

¹Throughout the paper we use \tilde{O} to hide poly-logarithmic factors in A_{tot} , t_{mix} , $1/(1-\gamma)$, $1/\epsilon$, and the number of states of the MDP.

We hope our method serves as a building block towards a more unified understanding the complexity of MDPs and RL. By providing a general SMD-based framework which is provably efficient for solving multiple prominent classes of MDPs we hope this paper may lead to a better understanding and broader application of the traditional convex optimization toolkit to modern RL. As a preliminary demonstration of flexibility of our framework, we show that it extends to yield new results for approximately optimizing constrained MDPs and hope it may find further utility.

1.1 Problem Setup

Throughout the paper we denote an MDP instance by a tuple $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathbf{r}, \gamma)$ with components defined as follows:

- \mathcal{S} - a finite set of states where each $i \in \mathcal{S}$ is called a *state of the MDP*, in tradition this is also denoted as s .
- $\mathcal{A} = \cup_{i \in [\mathcal{S}]} \mathcal{A}_i$ - a finite set of actions that is a collection of sets of actions \mathcal{A}_i for states $i \in \mathcal{S}$. We overload notation slightly and let $(i, a_i) \in \mathcal{A}$ denote *an action a_i at state i* . $A_{\text{tot}} := |\mathcal{A}| := \sum_{i \in \mathcal{S}} |\mathcal{A}_i|$ denotes the total number of state-action pairs.
- \mathcal{P} - the collection of state-to-state transition probabilities where $\mathcal{P} := \{p_{ij}(a_i) | i, j \in \mathcal{S}, a_i \in \mathcal{A}_i\}$ and $p_{ij}(a_i)$ denotes the probability of transition to state j when taking action a_i at state i .
- \mathbf{r} - the vector of state-action transitional rewards where $\mathbf{r} \in [0, 1]^{\mathcal{A}}$, r_{i, a_i} is the instant reward received when taking action a_i at state $i \in \mathcal{S}$.²
- γ - the discount factor of MDP, by which one down-weights the reward in the next future step. When $\gamma \in (0, 1)$, we call the instance a *discounted MDP* (DMDP) and when $\gamma = 1$, we call the instance an *average-reward MDP* (AMDP).

We use $\mathbf{P} \in \mathbb{R}^{\mathcal{A} \times \mathcal{S}}$ as the state-transition matrix where its (i, a_i) -th row corresponds to the transition probability from state $i \in \mathcal{S}$ where $a_i \in \mathcal{A}_i$ to state j . Correspondingly we use $\hat{\mathbf{I}}$ as the matrix with a_i -th row corresponding to \mathbf{e}_i , for all $i \in \mathcal{S}, a_i \in \mathcal{A}_i$.

Now, the model operates as follows: when at state i , one can pick an action a_i from the given action set \mathcal{A}_i . This generates a reward r_{i, a_i} . Also based on the transition model with probability $p_{ij}(a_i)$, it transits to state j and the process repeats. Our goal is to compute a random policy which determines which actions to take at each state. A random policy is a collection of probability distributions $\pi := \{\pi_i\}_{i \in \mathcal{S}}$, where $\pi_i \in \Delta^{\mathcal{A}_i}$ is a vector in the $|\mathcal{A}_i|$ -dimensional simplex

²The assumption that \mathbf{r} only depends on state action pair i, a_i is a common practice (Sidford et al., 2018a).

with $\pi_i(a_i)$ denoting the probability of taking $a_i \in \mathcal{A}_i$ at action j . One can extend π_i to the set of $\Delta^{\mathcal{A}}$ by filling in 0s on entries corresponding to other states $j \neq i$, and denote $\Pi \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ as the concatenated policy matrix with i -th row being the extended Δ_i . We denote \mathbf{P}^π as the transitional probability matrix of the MDP when using policy π , thus we have $\mathbf{P}^\pi(i, j) := \sum_{a_i \in \mathcal{A}_i} \pi_i(a_i) p_{ij}(a_i) = \Pi \cdot \mathbf{P}$ for all $i, j \in \mathcal{S}$, where \cdot in the right-hand side (RHS) denotes matrix-matrix multiplication. Further, we let \mathbf{r}^π denote corresponding average reward under policy π defined as $\mathbf{r}^\pi := \Pi \cdot \mathbf{r}$, where \cdot in RHS denotes matrix-vector multiplication. We overload notation and still use \mathbf{I} to denote the standard identity matrix if computing with regards to probability transition matrix $\Pi^\pi \in \mathbb{R}^{\mathcal{S} \times \mathcal{S}}$.

Given an MDP instance $\mathcal{M} := (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathbf{r}, \gamma)$ and an initial distribution over states $\mathbf{q} \in \Delta^{\mathcal{S}}$, we are interested in finding the optimal π^* among all policy π that maximizes the following cumulative reward \bar{v}^π of the MDP:

$$\pi^* := \arg \max_{\pi} \bar{v}^\pi \quad \text{where}$$

$$\bar{v}^\pi := \begin{cases} \mathbb{E}^\pi \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_{i_t, a_t} | i_1 \sim \mathbf{q} \right], & \text{i.e., DMDPs} \\ \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}^\pi \left[\sum_{t=1}^T r_{i_t, a_t} | i_1 \sim \mathbf{q} \right], & \text{i.e., AMDPs.} \end{cases}$$

Here $\{i_1, a_1, i_2, a_2, \dots, i_t, a_t\}$ are state-action transitions generated by the MDP under policy π . For the DMDP case, it also holds by definition that $\bar{v}^\pi := \mathbf{q}^\top (\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1} \mathbf{r}^\pi$.

For the AMDP case (i.e. when $\gamma = 1$), we also define $\boldsymbol{\nu}^\pi$ as the stationary distribution under policy π satisfying $\boldsymbol{\nu}^\pi = (\mathbf{P}^\pi)^\top \boldsymbol{\nu}^\pi$. To ensure the value of \bar{v}^π is well-defined, we restrict our attention to a subgroup which we call *mixing AMDP* satisfying the following mixing assumption:

Assumption A. An AMDP instance is mixing if t_{mix} , defined as follows, is bounded by $1/2$, i.e.

$$t_{\text{mix}} := \max_{\pi} \left[\arg \min_{t \geq 1} \max_{\mathbf{q} \in \Delta^{\mathcal{S}}} \|(\mathbf{P}^\pi)^\top \mathbf{q} - \boldsymbol{\nu}^\pi\|_1 \right] \leq \frac{1}{2}.$$

The mixing condition assumes for arbitrary policy π and arbitrary initial state, the resulting Markov chain leads toward a distribution close enough to its stationary distribution $\boldsymbol{\nu}^\pi$ starting from any initial state i in $O(t_{\text{mix}})$ time steps. This assumption implies the uniqueness of the stationary distribution, makes \bar{v}^π above well-defined with the equivalent $\bar{v}^\pi = (\boldsymbol{\nu}^\pi)^\top \mathbf{r}^\pi$, governing the complexity of our mixing AMDP algorithm, and is key for the results we prove (Theorem 1). Further, the assumption is equivalent as in Wang (2017b), up to constant factors.

By nature of the definition of *mixing AMDP*, we note that the value of a strategy π is independent of initial distribution \mathbf{q} and only dependent of the eventual stationary distribution

as long as the AMDP is mixing, which also implies \bar{v}^π is always well-defined. For this reason, sometimes we also omit $i_1 \sim \mathbf{q}$ in the corresponding definition of \bar{v}^π .

We call a policy π an ϵ -(*approximate*) *optimal policy* for the MDP problem, if it satisfies $\bar{v}^\pi \geq \bar{v}^* - \epsilon$.³ We call a policy an expected ϵ -(*approximate*) *optimal policy* if it satisfies the condition in expectation, i.e. $\mathbb{E}\bar{v}^\pi \geq \bar{v}^* - \epsilon$. The goal of paper is to develop efficient algorithms that find (expected) ϵ -approximate policy for the given MDP instance assuming access to a generative model.

1.2 Main Results

The main result of the paper is a unified framework based on randomized primal-dual stochastic mirror descent (SMD) that with high probability finds an (expected) ϵ -optimal policy with some sample complexity guarantee. Formally we provide two algorithms (see Algorithm 1 for both cases) with the following guarantees respectively.

Theorem 1. *Given a mixing AMDP tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathbf{r})$, let $\epsilon \in (0, 1)$, one can construct an expected ϵ -optimal policy π^ϵ from the decomposition (see Section 5) of output $\boldsymbol{\mu}^\epsilon$ of Algorithm 1 with sample complexity $O(t_{\text{mix}}^2 A_{\text{tot}} \epsilon^{-2} \log(A_{\text{tot}}))$.*

Theorem 2. *Given a DMDP tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathbf{r}, \gamma)$ with discount factor $\gamma \in (0, 1)$, let $\epsilon \in (0, 1)$, one can construct an expected ϵ -optimal policy π^ϵ from the decomposition (see Section 5) of output $\boldsymbol{\mu}^\epsilon$ of Algorithm 1 with sample complexity $O((1 - \gamma)^{-4} A_{\text{tot}} \epsilon^{-2} \log(A_{\text{tot}}))$.*

We remark that for both problems, the algorithm also gives with high probability an ϵ -optimal policy at the cost of an extra $\log(1/\delta)$ factor to the sample complexity through a reduction from high-probability to expected optimal policy (see Wang (2017b) for more details). Note that we only get randomized policies, and we leave the question of getting directly deterministic policies as an interesting open direction.

Table 1 gives a comparison of sample complexity between our methods and prior methods⁴ for computing an ϵ -approximate policy in DMDPs and AMDPs given a generative model.

As a generalization, we show how to solve constrained average-reward MDPs (cf. (Altman, 1999), a generalization of average-reward MDP) using the primal-dual stochastic mirror descent framework in Section 6. We build an algorithm that solves the constrained problem (13) to ϵ -accuracy within sample complexity

³Hereinafter, we use superscript $*$ and π^* interchangeably.

⁴Most methods assume a uniform action set \mathcal{A} for each of the $|\mathcal{S}|$ states, but can also be generalized to the non-uniform case parameterized by A_{tot} .

Algorithm 1 SMD for mixing AMDP / DMDPs

- 1: **Input:** MDP tuple $\mathcal{M} = (\mathcal{S}, \cup_{i \in \mathcal{S}} \mathcal{A}_i, \mathcal{P}, \mathbf{r}, \gamma)$, initial $(\mathbf{v}_0, \boldsymbol{\mu}_0) \in \mathbb{B}_{2M}^{\mathcal{S}} \times \Delta^{\mathcal{A}}$.
 - 2: **Output:** An expected ϵ -approximate solution $(\mathbf{v}^\epsilon, \boldsymbol{\mu}^\epsilon)$ for problem (5).
 - 3: **Parameter:** Step-size η^v, η^μ , number of iterations T , accuracy level ϵ .
 - 4: **for** $t = 1, \dots, T$ **do**
 - 5: // \mathbf{v} gradient estimation
 - 6: Sample $(i, a_i) \sim [\boldsymbol{\mu}]_{i, a_i}, j \sim p_{ij}(a_i), i' \sim q_{i'}$
 - 7: Set $\tilde{g}_{t-1}^v = \begin{cases} \mathbf{e}_j - \mathbf{e}_i & \text{mixing} \\ (1 - \gamma)\mathbf{e}_{i'} + \gamma\mathbf{e}_j - \mathbf{e}_i & \text{discounted} \end{cases}$
 - 8: // $\boldsymbol{\mu}$ gradient estimation
 - 9: Sample $(i, a_i) \sim \frac{1}{A_{\text{tot}}}, j \sim p_{ij}(a_i)$
 - 10: Set $\tilde{g}_{t-1}^\mu = \begin{cases} A_{\text{tot}}(v_i - v_j - r_{i, a_i})\mathbf{e}_{i, a_i} & \text{mixing} \\ A_{\text{tot}}(v_i - \gamma v_j - r_{i, a_i})\mathbf{e}_{i, a_i} & \text{discounted} \end{cases}$
 - 11: // Stochastic mirror descent steps (Π as projection)
 - 12: $\mathbf{v}_t \leftarrow \Pi_{\mathbb{B}_{2M}^{\mathcal{S}}}(\mathbf{v}_{t-1} - \eta^v \tilde{g}_{t-1}^v)$
 - 13: $\boldsymbol{\mu}_t \leftarrow \Pi_{\Delta^{\mathcal{A}}}(\boldsymbol{\mu}_{t-1} \circ \exp(-\eta^\mu \tilde{g}_{t-1}^\mu))$
 - 14: **end for**
 - 15: **Return** $(\mathbf{v}^\epsilon, \boldsymbol{\mu}^\epsilon) \leftarrow \frac{1}{T} \sum_{t \in [T]} (\mathbf{v}_t, \boldsymbol{\mu}_t)$
-

$O((t_{\text{mix}}^2 + K)D^2 A_{\text{tot}} \epsilon^{-2} \log(A_{\text{tot}}))$, where K and D^2 are number and size of the constraints. To the best of our knowledge this is the first sample complexity results for constrained MDPs given by a generative model.

As a byproduct, the framework we build in Section 3 also gives a stochastic algorithms (see Algorithm 2) that finds an expected ϵ -approximate solution of ℓ_∞ - ℓ_1 bilinear min-max problems of form $\min_{\mathbf{x} \in [-1, 1]^n} \max_{\mathbf{y} \in \Delta^m} \mathbf{y}^\top \mathbf{M} \mathbf{x} + \mathbf{b}^\top \mathbf{x} - \mathbf{c}^\top \mathbf{y}$ to ϵ -additive accuracy with runtime $\tilde{O}(((m + n) \|\mathbf{M}\|_\infty^2 + n \|\mathbf{b}\|_1^2 + m \|\mathbf{c}\|_\infty^2) \epsilon^{-2})$ given ℓ_1 sampler of iterate \mathbf{y} , and ℓ_1 samplers based on the input entries of \mathbf{M} , \mathbf{b} and \mathbf{c} (see Corollary 1 for details). Consequently, it solves (box constrained) ℓ_∞ regression problems of form $\min_{\mathbf{x} \in [-1, 1]^n} \|\mathbf{M} \mathbf{x} - \mathbf{c}\|_\infty$ to ϵ -additive accuracy within runtime $\tilde{O}(((m + n) \|\mathbf{M}\|_\infty^2 + m \|\mathbf{c}\|_\infty^2) \epsilon^{-2})$ given similar sampling access (see Remark 1 for details and Table 3 for comparison with previous results).

1.3 Technique Overview

We adopt the idea of formulating the MDP problem as a bilinear saddle point problem in light of linear duality, following the line of randomized model-free primal-dual π learning studied in Wang (2017a;b). This formulation relates MDP to solving bilinear saddle point problems with box and simplex domains, which falls into well-studied generalizations of convex optimization (Nemirovski, 2004a; Carmon et al., 2019).

Type	Method	Sample Complexity
mixing AMDP	Primal-Dual Method (Wang, 2017b)	$\tilde{O}(\tau^2 t_{\text{mix}}^2 A_{\text{tot}} \epsilon^{-2})$
	Our method (Theorem 1)	$\tilde{O}(t_{\text{mix}}^2 A_{\text{tot}} \epsilon^{-2})$
DMDP	Empirical QVI (Azar et al., 2012)	$\tilde{O}((1-\gamma)^{-5} A_{\text{tot}} \epsilon^{-2})$
	Primal-Dual Method (Wang, 2017a)	$\tilde{O}((1-\gamma)^{-6} \mathcal{S} ^2 A_{\text{tot}} \epsilon^{-2})$
	Primal-Dual Method (Wang, 2017a)	$\tilde{O}(\tau^4 (1-\gamma)^{-4} A_{\text{tot}} \epsilon^{-2})$
	Variance-reduced QVI (Sidford et al., 2018a)	$\tilde{O}((1-\gamma)^{-3} A_{\text{tot}} \epsilon^{-2})$
	Empirical MDP + Blackbox (Agarwal et al., 2020)	$\tilde{O}((1-\gamma)^{-3} A_{\text{tot}} \epsilon^{-2})$
	Variance-reduced Q-learning (Wainwright, 2019)	$\tilde{O}((1-\gamma)^{-3} A_{\text{tot}} \epsilon^{-2})$
	Our method (Theorem 2)	$\tilde{O}((1-\gamma)^{-4} A_{\text{tot}} \epsilon^{-2})$

Table 1. Comparison of sample complexity to get ϵ -optimal policy among stochastic methods (complete version see Appendix A). Here \mathcal{S} denotes state space, A_{tot} denotes number of state-action pair, t_{mix} is mixing time for mixing AMDP, and γ is discount factor for DMDP. Parameter τ shows up whenever the designed algorithm requires additional ergodic condition for MDP, i.e. there exists some distribution \mathbf{q} and $\tau > 0$ satisfying $\sqrt{1/\tau} \mathbf{q} \leq \nu^\pi \leq \sqrt{\tau} \mathbf{q}$, \forall policy π and its induced stationary distribution ν^π .

We study the efficiency of standard stochastic mirror descent (SMD) for this bilinear saddle point problem where the minimization (primal) variables are constrained to in a rescaled box domain and the maximization (dual) variables are constrained to lie in the simplex. We use the idea of local-norm variance bounds emerging in Shalev-Shwartz et al. (2012); Carmon et al. (2019), to design and analyze efficient stochastic estimators for the gradient of this problem that have low-variance under the corresponding local norms. We provide a new analytical way to bound the quality of an approximately-optimal policy constructed from the approximately optimal solution of bilinear saddle point problem, which utilizes the influence of the dual constraints under minimax optimality. Compared with prior work, we eliminate the assumption of ergodicity through extending the primal space by a constant size. Putting all these together, we prove a simple and natural SMD algorithm which solves both mixing AMDP and DMDP problems efficiently, with clear dependence on hardness parameters like mixing time t_{mix} or discount factor γ .

1.4 Related Work

1.4.1 ON SOLVING MDP

Within the tremendous body of study on MDPs, and more generally reinforcement learning, stands the well-studied classic problem of computational efficiency (i.e. iteration number, runtime, etc.) of finding optimal policy, given the entire MDP instance as an input. Traditional deterministic methods for the problems are value iteration, policy iteration, and linear programming. (Bertsekas et al., 1995; Ye, 2011), which find an approximately optimal policy to high-accuracy but have superlinear runtime in the usually high

problem dimension $\Omega(|\mathcal{S}| \cdot A_{\text{tot}})$.

To avoid the necessity of knowing the whole problem instance and having superlinear runtime dependence, more recently, researchers have designed stochastic algorithms assuming only a generative model that samples from state-transitions (Kakade et al., 2003). Azar et al. (2012) proved a lower bound of $\Omega((1-\gamma)^{-3} A_{\text{tot}} \epsilon^{-2})$ while also giving a Q-value-iteration algorithm with a higher guaranteed sample complexity. This was recently improved in Sidford et al. (2018b) using variance-reduction ideas, and was further improved to match (up to logarithmic factors) lower bound in (Sidford et al., 2018a) using a type of variance-reduced randomized value iteration. Soon later in Wainwright (2019), a variance-reduced Q-learning method also achieves nearly tight sample complexity for the discounted case. In Agarwal et al. (2020) the authors use a non-algorithmic approach that shows $\tilde{O}((1-\gamma)^{-3} A_{\text{tot}} \epsilon^{-2})$ samples suffice. On the other hand, Wang (2017a) designed a randomized primal-dual method, an instance of SMD with slightly different sampling distribution and form of update for the estimators, which has superlinear sample complexity guarantee unless additional ergodicity assumptions are made. Whether such randomized primal-dual methods necessarily incur this higher computational cost is unclear and a key motivation for our work.

While a few methods match (up to logarithmic factors) the lower bound shown for sample complexity for solving DMDP (Sidford et al., 2018a; Wainwright, 2019), it is unclear how and if one can design a similar method for average-reward MDP and get the optimal sample complexity dependence. The only related work under average-reward MDP setting uses primal-dual π -learning (Wang, 2017b), follow-

ing the stochastic primal-dual method in (Wang, 2017a). Being also a variant of SMD methods, their algorithm has a different domain setup, different update forms, and a more ad-hoc analysis compared with ours. The sample complexity of their method is $\tilde{O}(\tau^2 t_{\text{mix}}^2 A_{\text{tot}} \epsilon^{-2})$, which still requires an additional bound on the ergodicity parameter, $\tau > 0$, and depends on this parameter polynomially.

This mismatch of the theoretical efficiency of SMD methods as opposed to value iteration and Q-learning in solving DMDPs and the dependence of ergodicity when solving AMDPs with SMD methods motivated our study of a general SMD framework that could provably efficiently solve both problem instances. Table 1 includes a complete comparison between our results and the prior art for both cases.

1.4.2 ON ℓ_∞ REGRESSION AND BILINEAR SADDLE POINT PROBLEMS

Our framework gives a stochastic method for solving ℓ_∞ regression, which is a core problem in both combinatorics and continuous optimization due to its connection with maximum flow and linear programming (Lee & Sidford, 2014; 2015). Classic methods build on solving a smooth approximations of the problem (Nesterov, 2005) or finding the right regularizers and algorithms for its correspondingly primal-dual minimax problem (Nemirovski, 2004b; Nesterov, 2007). These methods have recently been improved to $\tilde{O}(\text{nnz} \|\mathbf{M}\|_\infty \epsilon^{-1})$ using a joint regularizer with nice area-convexity properties in Sherman (2017) or using accelerated coordinate method with a matching runtime bound in sparse-column case in Sidford & Tian (2018).

In comparison to all the state-of-the-art, for dense input data matrix our method gives the first algorithm with sublinear runtime dependence $O(m+n)$ instead of $O(\text{nnz})$. For completeness we include a comparison of runtimes for methods mentioned above in Appendix B (see Table 3).

Our sublinear method for ℓ_∞ -regression is closely related to a line of work on obtaining efficient stochastic methods for approximately solving *matrix games*, i.e. bilinear saddle point problems (Grigoriadis & Khachiyan, 1995; Clarkson et al., 2012; Palaniappan & Bach, 2016), and, in particular, a recent line of work by the authors and collaborators (Carmon et al., 2019; 2020) that explores the benefit of careful sampling and variance reduction in matrix games. In Carmon et al. (2019) we provide a framework to analyze variance-reduced SMD under local norms to obtain better complexity bounds for different domain setups, i.e. ℓ_1 - ℓ_1 , ℓ_1 - ℓ_2 , and ℓ_2 - ℓ_2 where ℓ_1 corresponds to the simplex and ℓ_2 corresponds to the Euclidean ball. In Carmon et al. (2020) we study the improved sublinear and variance-reduced coordinate methods for these domain setups utilizing the design of optimal gradient estimators. This paper adapts the local norm analysis and coordinate-wise gradient estimator

design in Carmon et al. (2019; 2020) to obtain our SMD algorithm and analysis for ℓ_1 - ℓ_∞ games.

2 Preliminaries

First, we introduce several known tools for studying MDPs.

2.1 Bellman Equation.

For mixing AMDP, \bar{v}^* is the optimal average reward if and only if there exists a vector $\mathbf{v}^* = (v_i^*)_{i \in \mathcal{S}}$ satisfying its corresponding *Bellman equation* (Bertsekas et al., 1995)

$$\bar{v}^* + v_i^* = \max_{a_i \in \mathcal{A}_i} \left\{ \sum_{j \in \mathcal{S}} p_{ij}(a_i) v_j^* + r_{i,a_i} \right\}, \forall i \in \mathcal{S}. \quad (1)$$

When considering a mixing AMDP as in the paper, the existence of solution to the above equation can be guaranteed. However, it is important to note that one cannot guarantee the uniqueness of the optimal \mathbf{v}^* . In fact, for each optimal solution \mathbf{v}^* , $\mathbf{v}^* + c\mathbf{1}$ is also an optimal solution.

For DMDP, one can show that at optimal policy π^* , each state $i \in \mathcal{S}$ can be assigned an optimal cost-to-go value v_i^* satisfying the following *Bellman equation* (Bertsekas et al., 1995)

$$v_i^* = \max_{a_i \in \mathcal{A}_i} \left\{ \sum_{j \in \mathcal{S}} \gamma p_{ij}(a_i) v_j^* + r_{i,a_i} \right\}, \forall i \in \mathcal{S}. \quad (2)$$

When $\gamma \in (0, 1)$, it is straightforward to guarantee the existence and uniqueness of the optimal solution $\mathbf{v}^* := (v_i^*)_{i \in \mathcal{S}}$ to the system.

2.2 Linear Programming (LP) Formulation.

We can further write the above Bellman equations equivalently as the following primal or dual linear programming problems. We define the domain as $\mathbb{B}_m^S := m \cdot [-1, 1]^S$ where \mathbb{B} stands for box, and $\Delta^n := \{\Delta \in \mathbb{R}^n, \Delta_i \geq 0, \sum_{i \in [n]} \Delta_i = 1\}$ for standard n -dimension simplex.

For mixing AMDP case, the linear programming formulation leveraging matrix notation is (with (P), (D) representing (equivalently) the primal form and the dual form respectively)

$$\begin{aligned} \text{(P)} \quad & \min_{\bar{v}, \mathbf{v}} && \bar{v} \\ & \text{subject to} && \bar{v} \cdot \mathbf{1} + (\hat{\mathbf{I}} - \mathbf{P})\mathbf{v} - \mathbf{r} \geq 0, \\ \text{(D)} \quad & \max_{\boldsymbol{\mu} \in \Delta^{\mathcal{A}}} && \boldsymbol{\mu}^\top \mathbf{r} \\ & \text{subject to} && (\hat{\mathbf{I}} - \mathbf{P})^\top \boldsymbol{\mu} = \mathbf{0}. \end{aligned} \quad (3)$$

The optimal values of both systems are the optimal expected cumulative reward \bar{v}^* under optimal policy π^* , thus here-

inafter we use \bar{v}^* and \bar{v}^{π^*} interchangeably. Given the optimal dual solution μ^* , one can without loss of generality impose the constraint of $\langle \mathbf{I}^\top \mu^*, \mathbf{v}^* \rangle = 0$ ⁵ to ensure uniqueness of the primal problem (P).

For DMDP case, the equivalent linear programming is

$$\begin{aligned}
 \text{(P)} \quad & \min_{\mathbf{v} \in \mathbb{B}_{2M}^S} (1 - \gamma) \mathbf{q}^\top \mathbf{v} \\
 & \text{subject to} \quad (\hat{\mathbf{I}} - \gamma \mathbf{P}) \mathbf{v} - \mathbf{r} \geq 0, \\
 \text{(D)} \quad & \max_{\mu \in \Delta^A} \mu^\top \mathbf{r} \\
 & \text{subject to} \quad (\hat{\mathbf{I}} - \gamma \mathbf{P})^\top \mu = (1 - \gamma) \mathbf{q}.
 \end{aligned} \tag{4}$$

Given a fixed initial distribution \mathbf{q} , the optimal values of both systems are a $(1 - \gamma)$ factor of the optimal expected cumulative reward, i.e. $(1 - \gamma) \bar{v}^*$ under optimal policy π^* .

2.3 Minimax Formulation.

By standard linear duality, we can recast the problem formulation in Section 2.2 using the method of Lagrangian multipliers, as bilinear saddle-point (minimax) problems. For AMDPs the minimax formulation is

$$\begin{aligned}
 \min_{\bar{v}, \mathbf{v} \in \mathbb{B}_{2M}^S} \max_{\mu \in \Delta^A} f(\bar{v}, \mathbf{v}, \mu), \\
 \text{where } f(\bar{v}, \mathbf{v}, \mu) & := \bar{v} + \mu^\top (-\bar{v} \cdot \mathbf{1} + (\mathbf{P} - \hat{\mathbf{I}}) \mathbf{v} + \mathbf{r}) \\
 & = \mu^\top ((\mathbf{P} - \hat{\mathbf{I}}) \mathbf{v} + \mathbf{r})
 \end{aligned} \tag{5}$$

For DMDPs the minimax formulation is

$$\begin{aligned}
 \min_{\mathbf{v} \in \mathbb{B}_{2M}^S} \max_{\mu \in \Delta^A} f_{\mathbf{q}}(\mathbf{v}, \mu), \\
 \text{where } f_{\mathbf{q}}(\mathbf{v}, \mu) & := (1 - \gamma) \mathbf{q}^\top \mathbf{v} + \mu^\top ((\gamma \mathbf{P} - \hat{\mathbf{I}}) \mathbf{v} + \mathbf{r}).
 \end{aligned} \tag{6}$$

Note in both cases we have added the constraint of $\mathbf{v} \in \mathbb{B}_{2M}^S$. The M is different for each case, and will be specified in Section 4 to ensure that $\mathbf{v}^* \in \mathbb{B}_{2M}^S$. As a result, constraining the bilinear saddle point problem on a restricted domain for primal variables will not affect the optimality of the original optimal solution due to its global optimality, but will considerably save work for the algorithm by considering a smaller domain.

For each problem we define the duality gap of the minimax problem $\min_{\mathbf{v} \in \mathbb{B}_{2M}^S} \max_{\mu \in \Delta^A} f(\mathbf{v}, \mu)$ at a given pair of feasible solution (\mathbf{v}, μ) as $\text{Gap}(\mathbf{v}, \mu) := \max_{\mu' \in \Delta^A} f(\mathbf{v}, \mu') - \min_{\mathbf{v}' \in \mathbb{B}_{2M}^S} f(\mathbf{v}', \mu)$.

An ϵ -approximate solution of the minimax problem is a pair of feasible solution $(\mathbf{v}^\epsilon, \mu^\epsilon) \in \mathbb{B}_{2M}^S \times \Delta^A$ with its duality gap bounded by ϵ , i.e. $\text{Gap}(\mathbf{v}^\epsilon, \mu^\epsilon) \leq \epsilon$. An expected ϵ -approximate solution is one satisfying $\mathbb{E} \text{Gap}(\mathbf{v}^\epsilon, \mu^\epsilon) \leq \epsilon$.

⁵ $\hat{\mathbf{I}}^\top \mu^*$ represents the stationary distribution over states given optimal policy π^* constructed from optimal dual variable μ^* .

3 Stochastic Mirror Descent Framework

In this section, we consider the following ℓ_∞ - ℓ_1 bilinear games as an abstraction of the MDP minimax problems of interest. Such games are induced by one player minimizing over the box domain (ℓ_∞) and the other maximizing over the simplex domain (ℓ_1) a bilinear objective:

$$\min_{\mathbf{x} \in \mathbb{B}_b^n} \max_{\mathbf{y} \in \Delta^m} f(\mathbf{x}, \mathbf{y}) := \mathbf{y}^\top \mathbf{M} \mathbf{x} + \mathbf{b}^\top \mathbf{x} - \mathbf{c}^\top \mathbf{y}, \tag{7}$$

We study the efficiency of coordinate stochastic mirror descent algorithms onto this ℓ_∞ - ℓ_1 minimax problem. The analysis follows from extending a fine-grained analysis of mirror descent with Bregman divergence using local norm arguments in Shalev-Shwartz et al. (2012); Carmon et al. (2019) to the ℓ_∞ - ℓ_1 domain. We defer all proofs in this section to Appendix C.

At a given iterate $(\mathbf{x}, \mathbf{y}) \in B_b^n \times \Delta^m$, our algorithm computes an estimate of the gradients for both sides defined as $g^x(\mathbf{x}, \mathbf{y}) := \mathbf{M}^\top \mathbf{y} + \mathbf{b} \in \mathbb{R}^n$ (x -gradient, or g^x); $g^y(\mathbf{x}, \mathbf{y}) := -\mathbf{M} \mathbf{x} + \mathbf{c} \in \mathbb{R}^m$ (y -gradient, or g^y).

The norm we use to measure these gradients are induced by Bregman divergence, a natural extension of Euclidean norm. For our analysis we choose to use $V_x(\mathbf{x}') := \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2$ and $V_y(\mathbf{y}') := \sum_{i \in [m]} y_i \log(\frac{y'_i}{y_i})$ (KL-divergence), which are also common practice (Wang, 2017a,b; Nesterov, 2005) catering to the geometry of each domain, and induce the dual norms on the gradients in form $\|g^x\| := \|g^x\|_2 = \sqrt{\sum_{j \in [n]} g_j^{x2}}$ (standard ℓ_2 -norm) for x side, and $\|g^y\|_{\mathbf{y}'} := \sum_{i \in [m]} y'_i (g_i^y)^2$ (a weighted ℓ_2 -norm) for y side.

To describe the properties of estimators needed for our algorithm, we introduce the following definition of *bounded estimator* as follows.

Definition 1 (Bounded Estimator). *Given the following properties on mean, scale and variance of an estimator:*

- (i) *unbiasedness:* $\mathbb{E} \tilde{g} = g$;
 - (ii) *bounded maximum entry:* $\|\tilde{g}\|_\infty \leq c$ with probability 1;
 - (iii) *bounded second-moment:* $\mathbb{E} \|\tilde{g}\|^2 \leq v$
- we call \tilde{g} a $(c, v, \|\cdot\|)$ -bounded estimator if satisfying (i) and (iii), call it and a $(c, v, \|\cdot\|_\Delta^m)$ -bounded estimator if besides (i) and (ii), it also satisfies (iii) with local norm $\|\cdot\|_{\mathbf{y}}$ for all $\mathbf{y} \in \Delta^m$.*

Now we give Algorithm 2, our general algorithmic framework for solving (7) given efficient bounded estimators for the gradient. Its theoretical guarantees are given in Theorem 3 which bounds the number of iterations needed to obtain expected ϵ -approximate solution.

Theorem 3. *Given an ℓ_∞ - ℓ_1 game, i.e. (7), and desired accuracy ϵ , a $(v^x, \|\cdot\|_2)$ -bounded estimator \tilde{g}^x , and a $(\frac{4v^y}{\epsilon}, v^y, \|\cdot\|_{\Delta^m})$ -bounded estimator \tilde{g}^y , Algorithm 2 with*

Algorithm 2 SMD for ℓ_∞ - ℓ_1 game

- 1: **Input:** Desired accuracy ϵ , primal domain size b , $(v^x, \|\cdot\|_2)$ -bounded estimator g^x , $(\frac{20v^y}{\epsilon}, v^y, \|\cdot\|_{\Delta^m})$ -bounded estimator g^y
- 2: **Output:** An expected ϵ -approximate solution $(\mathbf{x}^\epsilon, \mathbf{y}^\epsilon)$ for problem (7).
- 3: **Parameter:** Step-size η^x, η^y , total iteration number T .
- 4: **for** $t = 0, \dots, T - 1$ **do**
- 5: Get \tilde{g}^x estimator for g^x , \tilde{g}^y estimator for g^y
- 6: Update $\mathbf{x}_{t+1} \leftarrow \arg \min_{\mathbf{x} \in \mathbb{B}_b^n} \langle \eta^x \tilde{g}^x(\mathbf{x}_t, \mathbf{y}_t), \mathbf{x} \rangle + V_{\mathbf{x}_t}(\mathbf{x})$
- 7: Update $\mathbf{y}_{t+1} \leftarrow \arg \min_{\mathbf{y} \in \Delta^m} \langle \eta^y \tilde{g}^y(\mathbf{x}_t, \mathbf{y}_t), \mathbf{y} \rangle + V_{\mathbf{y}_t}(\mathbf{y})$
- 8: **end for**
- 9: **Return** $(\mathbf{x}^\epsilon, \mathbf{y}^\epsilon) \leftarrow \frac{1}{T} \sum_{t \in [T]} (\mathbf{x}_t, \mathbf{y}_t)$

choice of parameters $\eta_x \leq \frac{\epsilon}{4v^x}$, $\eta_y \leq \frac{\epsilon}{4v^y}$ outputs an expected ϵ -approximate optimal solution within any iteration number $T \geq \max\{\frac{16nb^2}{\epsilon\eta_x}, \frac{8\log m}{\epsilon\eta_y}\}$.

Sketch of Proof For simplicity we only include a proof sketch here and defer the complete proof details to Appendix C.2.

Regret bounds with local norms. The core statement is a standard regret bound using local norms (see Lemma 10 and Lemma 11) which summing together gives the following guarantee (let $\tilde{g}_t^x, \tilde{g}_t^y$ denote $\tilde{g}^x(\mathbf{x}_t, \mathbf{y}_t), \tilde{g}^y(\mathbf{x}_t, \mathbf{y}_t)$)

$$\begin{aligned} & \sum_{t \in [T]} \langle \tilde{g}_t^x, \mathbf{x}_t - \mathbf{x} \rangle + \sum_{t \in [T]} \langle \tilde{g}_t^y, \mathbf{y}_t - \mathbf{y} \rangle \\ & \leq \frac{V_{\mathbf{x}_0}(\mathbf{x})}{\eta^x} + \frac{\sum_{t=0}^T \eta^x \|\tilde{g}_t^x\|_2^2}{2} + \frac{V_{\mathbf{y}_0}(\mathbf{y})}{\eta^y} + \frac{\sum_{t=0}^T \eta^y \|\tilde{g}_t^y\|_{\mathbf{y}_t}^2}{2}. \end{aligned} \quad (8)$$

Note one needs the bounded maximum entry condition for \tilde{g}^y as the condition to use Lemma 11.

Domain size. The domain size can be bounded as $\max_{\mathbf{x} \in \mathbb{B}_b^n} V_{\mathbf{x}_0}(\mathbf{x}) \leq 2nb^2$, $\max_{\mathbf{y} \in \Delta^m} V_{\mathbf{y}_0}(\mathbf{y}) \leq \log m$ by definition of their corresponding Bregman divergences.

Second-moment bounds. This is given through the bounded second-moment properties of estimators directly.

*Ghost-iterate analysis.*⁶ In order to substitute \tilde{g}^x, \tilde{g}^y with g^x, g^y for LHS of Eq. (8), one can apply the regret bounds again to ghost iterates generated by taking gradient step with $\hat{g} = g - \tilde{g}$ coupled with each iteration. The additional terms coming from this extra regret bounds are in expectation 0 through conditional expectation computation.

Optimal tradeoff. One pick η_x, η_y, T accordingly to get the desired guarantee as stated in Theorem 3. \square

⁶For standard SMD on convex problems this step is unnecessary. One can directly use conditional expectation by fixing $\mathbf{x} = \mathbf{x}^*$. However, for saddle-point problems, the same technique only gives $\max_{\mathbf{x}, \mathbf{y}} \mathbb{E}[\text{regret}] \leq \epsilon$. The ghost iterates analysis is standard (Nemirovski, 2004a; Carmon et al., 2019) and necessary to get a bound in terms of $\mathbb{E} \max_{\mathbf{x}, \mathbf{y}} [\text{regret}] \leq \epsilon$ instead.

Now we design gradient estimators assuming certain sampling oracles to ensure good bounded properties.

When $\mathbf{x} \in \mathbb{B}_1^n$, this leads to the theoretical guarantee as stated formally in Corollary 1. We defer design of estimators and all proofs to Appendix C.3 and simply state the theoretical runtime guarantee of Algorithm 2 here.

Corollary 1. *Given an ℓ_∞ - ℓ_1 game (7) with domains $\mathbf{x} \in \mathbb{B}_1^n, \mathbf{y} \in \Delta^m, \epsilon \in (0, 1)$ and $\|\mathbf{M}\|_\infty + \|\mathbf{c}\|_\infty = \Omega(1)$. If one has all sampling oracles needed⁷, Algorithm 2 with certain gradient estimators (see (19) and (20)) finds an expected ϵ -approximate solution in runtime (equivalent as sample complexity here) $O((n + m \log m) \|\mathbf{M}\|_\infty^2 + n \|\mathbf{b}\|_1^2 + m \log m \|\mathbf{c}\|_\infty^2) \cdot \epsilon^{-2}$.*

Finally, we remark that one can also use Algorithm 2 to solve ℓ_∞ -regression, i.e. the problem of finding $\mathbf{x}^* := \arg \min_{\mathbf{x} \in \mathbb{B}_1^n} \|\mathbf{M}\mathbf{x} - \mathbf{c}\|_\infty$ by simply writing it in equivalent minimax form of $\min_{\mathbf{x} \in \mathbb{B}_1^n} \max_{\mathbf{y} \in \Delta^m} \mathbf{y}^\top (\hat{\mathbf{M}}\mathbf{x} - \hat{\mathbf{c}})$ where $\hat{\mathbf{M}} := [\mathbf{M}; -\mathbf{M}]$ and $\hat{\mathbf{c}} := [\mathbf{c}; -\mathbf{c}]$.

Remark 1. *Algorithm 2 produces an expected ϵ -approximate solution \mathbf{x}^ϵ satisfying $\mathbb{E} \|\mathbf{M}\mathbf{x}^\epsilon - \mathbf{c}\|_\infty \leq \|\mathbf{M}\mathbf{x}^* - \mathbf{c}\|_\infty + \epsilon$, within runtime $\tilde{O}\left(\left[(m+n)\|\mathbf{M}\|_\infty^2 + m\|\mathbf{c}\|_\infty^2\right] \cdot \epsilon^{-2}\right)$.*

4 Gradient Estimators for Mixing AMDP and DMDP

Specifically, both MDP problems induced by solving mixing or DMDP are in minimax form of (7) if we let $\mathbf{y} \leftarrow \boldsymbol{\mu}$ with $m = A_{\text{tot}}, \mathbf{x} \leftarrow \mathbf{v}$ with $n = |\mathcal{S}|$ and $b \leftarrow 2M$ with M chosen to be $M = 3t_{\text{mix}}$ for mixing AMDP and $M = \frac{1}{1-\gamma}$ for DMDP so that $\mathbf{v}^* \in \mathbb{B}_M^S$; we defer readers to proofs of Lemma 6 and 9 for a complete argument on it.

In this section, we give a cleaner way to construct gradient estimators with desired properties for mixing and discounted cases utilizing problem structure and the generative model at hand. Such a gradient estimator samples state-action pair for \mathbf{v} -side using a dynamic distribution induced by $\boldsymbol{\mu}$, while sampling state-action pair for $\boldsymbol{\mu}$ -side using a uniform distribution. We defer all proofs in this section to Appendix D.

4.1 Mixing AMDPs

For the mixing case, we set $M = 3t_{\text{mix}}$ to guarantee $\mathbf{v}^* \in \mathbb{B}_{2M}^S$.⁸ This follows immediately from a lemma that relates

⁷Note all the sampling oracles needed are essentially ℓ_1 samplers proportional to the matrix / vector entries, and an ℓ_1 sampler induced by $y \in \Delta^m$.

⁸Note the one can show $\|\mathbf{v}^*\|_\infty \leq M$ and the extra coefficient 2 in box size $2M$ is to ensure a stronger condition on dual constraints for approximate solutions, which can be seen more clearly in proof of Lemma 6

matrix norm of interest to the mixing property of MDP, which we defer readers to Appendix D.1 for details.

Given domain setups, now we describe formally the gradient estimators used in Algorithm 1 and their properties.

For the \mathbf{v} -side, we consider the following gradient estimator

$$\begin{aligned} \text{Sample } (i, a_i) &\sim [\boldsymbol{\mu}]_{i, a_i}, j \sim p_{ij}(a_i). \\ \text{Set } \tilde{g}^{\mathbf{v}}(\mathbf{v}, \boldsymbol{\mu}) &= \mathbf{e}_j - \mathbf{e}_i. \end{aligned} \quad (9)$$

This is a bounded gradient estimator for the box domain.

Lemma 1. $\tilde{g}^{\mathbf{v}}$ as in (9) is a $(2, \|\cdot\|_2)$ -bounded estimator.

For the $\boldsymbol{\mu}$ -side, we consider the following gradient estimator

$$\begin{aligned} \text{Sample } (i, a_i) &\sim 1/A_{\text{tot}}, j \sim p_{ij}(a_i). \\ \text{Set } \tilde{g}^{\boldsymbol{\mu}}(\mathbf{v}, \boldsymbol{\mu}) &= A_{\text{tot}}(v_i - \gamma v_j - r_{i, a_i})\mathbf{e}_{i, a_i}. \end{aligned} \quad (10)$$

This is a bounded gradient estimator for the simplex domain.

Lemma 2. $\tilde{g}^{\boldsymbol{\mu}}$ defined in (10) is a $((2M+1)A_{\text{tot}}, 9(M^2+1)A_{\text{tot}}, \|\cdot\|_{\Delta^{\mathcal{A}}})$ -bounded estimator.

Theorem 3 together with guarantees of designed gradient estimators in Lemma 1, 2 and choice of $M = 3t_{\text{mix}}$ gives Corollary 4.

Corollary 2. Given mixing AMDP tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathbf{r})$ with desired accuracy $\epsilon \in (0, 1)$, Algorithm 1 with parameter choice $\eta^{\mathbf{v}} = O(\epsilon)$, $\eta^{\boldsymbol{\mu}} = O(\epsilon t_{\text{mix}}^{-2} A_{\text{tot}}^{-1})$ outputs an expected ϵ -approximate solution to mixing minimax problem (5) with sample complexity $O(t_{\text{mix}}^2 A_{\text{tot}} \epsilon^{-2} \log(A_{\text{tot}}))$.

The proof follows immediately by noticing each iteration costs $O(1)$ sample generation, thus directly transferring the total iteration number to sample complexity.

4.2 DMDPs

Here, we pick $M = (1 - \gamma)^{-1}$, with the guarantee that $\mathbf{v}^* \in \mathbb{B}_{2M}^{\mathcal{S}}$ following from Lemma 15 we state and prove in Appendix D.1.

For discounted case one construct gradient estimators in a similar way. For the \mathbf{v} -side, we consider the following gradient estimator

$$\begin{aligned} \text{Sample } (i, a_i) &\sim [\boldsymbol{\mu}]_{i, a_i}, j \sim p_{ij}(a_i), i' \sim q_{i'} \\ \text{Set } \tilde{g}^{\mathbf{v}}(\mathbf{v}, \boldsymbol{\mu}) &= (1 - \gamma)\mathbf{e}_{i'} + \gamma\mathbf{e}_j - \mathbf{e}_i. \end{aligned} \quad (11)$$

Lemma 3. $\tilde{g}^{\mathbf{v}}$ as in (11) is a $(2, \|\cdot\|_2)$ -bounded estimator.

For the $\boldsymbol{\mu}$ -side, we consider the following gradient estimator

$$\begin{aligned} \text{Sample } (i, a_i) &\sim \frac{1}{A_{\text{tot}}}, j \sim p_{ij}(a_i). \\ \text{Set } \tilde{g}^{\boldsymbol{\mu}}(\mathbf{v}, \boldsymbol{\mu}) &= A_{\text{tot}}(v_i - \gamma v_j - r_{i, a_i})\mathbf{e}_{i, a_i}. \end{aligned} \quad (12)$$

Lemma 4. $\tilde{g}^{\boldsymbol{\mu}}$ defined in (12) is a $((2M+1)A_{\text{tot}}, 9(M^2+1)A_{\text{tot}}, \|\cdot\|_{\Delta^{\mathcal{A}}})$ -bounded estimator.

Theorem 3 together with guarantees of gradient estimators in use in Lemma 3, 4 and choice of $M = (1 - \gamma)^{-1}$ gives Corollary 3.

Corollary 3. Given DMDP tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathbf{r}, \gamma)$ with desired accuracy $\epsilon \in (0, 1)$, Algorithm 1 outputs an expected ϵ -approximate solution to discounted minimax problem (6) with sample complexity $O((1 - \gamma)^{-2} A_{\text{tot}} \epsilon^{-2} \log(A_{\text{tot}}))$.

5 From Optimal Solution to Optimal Policy

In this section, we relate the quality of an approximate solution to minimax problem to the quality of approximate policy one can construct from it, for both cases. To do that, we show that the enlarged primal space (size $2M$ instead of M) allows one to bound the quality of dual variables $\boldsymbol{\mu}^\epsilon$ better. We defer all proofs of this section to Appendix E.

5.1 Mixing AMDPs

Now we proceed to show how to convert an ϵ -approximate solution of (5) to an $\Theta(\epsilon)$ -approximate policy for (3).

First we introduce a lemma that relates the dual variable $\boldsymbol{\mu}^\epsilon$ with optimal cost-to-go values \mathbf{v}^* and expected reward \bar{v}^* .

Lemma 5. If $(\mathbf{v}^\epsilon, \boldsymbol{\mu}^\epsilon)$ is an expected ϵ -approximate optimal solution to mixing AMDP minimax problem (5), then for any optimal \mathbf{v}^* and \bar{v}^* , $\mathbb{E} \left[\boldsymbol{\mu}^{\epsilon \top} \left[(\hat{\mathbf{I}} - \mathbf{P})\mathbf{v}^* - \mathbf{r} \right] + \bar{v}^* \right] \leq \epsilon$.

Next we transfer an optimal solution to an optimal policy, formally through Lemma 6.

Lemma 6. Given an ϵ -approximate solution $(\mathbf{v}^\epsilon, \boldsymbol{\mu}^\epsilon)$ for mixing minimax problem as defined in (5), let π^ϵ be the unique decomposition (in terms of $\boldsymbol{\lambda}^\epsilon$) such that $\mu_{i, a_i}^\epsilon = \lambda_i^\epsilon \cdot \pi_{i, a_i}^\epsilon, \forall i \in \mathcal{S}, a_i \in \mathcal{A}_i$, where $\boldsymbol{\lambda} \in \Delta^{\mathcal{S}}, \pi_i \in \Delta^{\mathcal{A}_i}, \forall i \in \mathcal{S}$. Taking $\pi := \pi^\epsilon$ as our policy, it holds that $\bar{v}^* \leq \mathbb{E} \bar{v}^\pi + 3\epsilon$.

To prove the lemma, we need the following helper lemma bounding the matrix norm.

Lemma 7. Given a mixing AMDP, policy π , and its probability transition matrix $\mathbf{P}^\pi \in \mathbb{R}^{\mathcal{S} \times \mathcal{S}}$ and stationary distribution $\boldsymbol{\nu}^\pi$, $\|(\mathbf{I} - \mathbf{P}^\pi + \mathbf{1}(\boldsymbol{\nu}^\pi)^\top)^{-1}\|_\infty \leq 3t_{\text{mix}}$.

Using this fact one can prove Lemma 6 by showing the linear constraints in dual formulation (D) of (3) are approximately satisfied given an ϵ -approximate optimal solution $(\mathbf{v}^\epsilon, \boldsymbol{\mu}^\epsilon)$ to minimax problem (5).

Lemma 6 shows one can construct an expected ϵ -optimal policy from an expected $\epsilon/3$ -approximate solution of the minimax problem (5). Thus, using Corollary 4 one directly obtains our desired total sample complexity for Algorithm 1

to solve mixing AMDPs to desired accuracy, as stated in Theorem 1.

5.2 DMDPs

Now we show how to convert an ϵ -approximate solution of (6) to an $\Theta((1 - \gamma)^{-1}\epsilon)$ -approximate policy of (4).

First we introduce a lemma similar to Lemma 5 that relates the dual variable $\boldsymbol{\mu}^\epsilon$ with values \mathbf{v}^* under ϵ -approximation.

Lemma 8. *If $(\mathbf{v}^\epsilon, \boldsymbol{\mu}^\epsilon)$ is an ϵ -approximate optimal solution to the DMDP minimax problem (6), then for optimal \mathbf{v}^* , $\mathbb{E}\boldsymbol{\mu}^{\epsilon\top} \left[(\hat{\mathbf{I}} - \gamma\mathbf{P})\mathbf{v}^* - \mathbf{r} \right] \leq \epsilon$.*

Next we transfer an optimal solution to an optimal policy, formally through Lemma 9.

Lemma 9. *Given an expected ϵ -approximate solution $(\mathbf{v}^\epsilon, \boldsymbol{\mu}^\epsilon)$ for discounted minimax problem as defined in (6), let π^ϵ be the unique decomposition (in terms of $\boldsymbol{\lambda}^\epsilon$) such that $\mu_{i,a_i}^\epsilon = \lambda_i^\epsilon \cdot \pi_{i,a_i}^\epsilon, \forall i \in \mathcal{S}, a_i \in \mathcal{A}_i$, where $\boldsymbol{\lambda} \in \Delta^{\mathcal{S}}, \pi_i^\epsilon \in \Delta^{\mathcal{A}_i}, \forall i \in \mathcal{S}$. Taking $\pi := \pi^\epsilon$ as our policy, it holds that $\bar{v}^* \leq \mathbb{E}\bar{v}^\pi + 3\epsilon/(1 - \gamma)$.*

Lemma 9 shows it suffices to find an expected $(1 - \gamma)\epsilon$ -approximate solution to problem (6) to get an expected ϵ -optimal policy. Together with Corollary 3 this directly yields the sample complexity as claimed in Theorem 2.

6 Constrained MDP

In this section, we consider solving a generalization of the mixing AMDP problem with additional linear constraints, which has been an important and well-known problem class along the study of MDP (Altman, 1999); we defer readers to Appendix G for derivation omitted in this section.

Formally, we focus on approximately solving the following dual formulation of constrained mixing AMDPs⁹:

$$\begin{aligned} \text{(D)} \quad & \max_{\boldsymbol{\mu} \in \Delta^{\mathcal{A}}} && 0 \\ & \text{subject to} && (\hat{\mathbf{I}} - \mathbf{P})^\top \boldsymbol{\mu} = \mathbf{0}, \quad \mathbf{D}^\top \boldsymbol{\mu} \geq \mathbf{1}, \end{aligned} \quad (13)$$

where $\mathbf{D} = [\mathbf{d}_1 \ \cdots \ \mathbf{d}_K]$ under the additional assumptions that $\mathbf{d}_k \geq \mathbf{0}, \forall k \in [K]$ and the problem is strictly feasible (with an inner point in its feasible set). Our goal is to compute ϵ -approximate policies and solutions for (13) defined as follows.

Definition 2. *Given a policy π with its stationary distribution $\boldsymbol{\nu}^\pi$, it is an ϵ -approximate policy of system (13) if for $\boldsymbol{\mu}$ defined as $\mu_{i,a_i} = \nu_i^\pi \pi_{i,a_i}, \forall i \in \mathcal{S}, a_i \in \mathcal{A}_i$ it is an*

⁹One can reduce the general case of $\mathbf{D}^\top \boldsymbol{\mu} \geq \mathbf{c}$ for some $\mathbf{c} > \mathbf{0}$ to this case by taking $\mathbf{d}_k \leftarrow \mathbf{d}_k / c_k$, under which an ϵ -approximate solution as defined in (14) of the modified problem corresponds to a multiplicatively approximate solution satisfying $\mathbf{D}^\top \boldsymbol{\mu} \geq (1 - \epsilon)\mathbf{c}$.

ϵ -approximate solution of (13), i.e. it satisfies

$$\boldsymbol{\mu}^\top (\hat{\mathbf{I}} - \mathbf{P}) = \mathbf{0}, \quad \mathbf{D}^\top \boldsymbol{\mu} \geq (1 - \epsilon)\mathbf{1}. \quad (14)$$

For $D := \|\mathbf{D}\|_{\max} := \max_{i,a_i,k} |[d_k]_{i,a_i}|$ and $M := 3Dt_{\text{mix}}$ we consider the following equivalent problem:

$$\min_{\mathbf{v} \in \mathbb{B}_{2M}^{\mathcal{S}}, \mathbf{s}: \sum_k s_k \leq 2, \mathbf{s} \geq \mathbf{0}} \max_{\boldsymbol{\mu} \in \Delta^{\mathcal{A}}} f(\mathbf{v}, \mathbf{s}, \boldsymbol{\mu}) \quad (15)$$

$$\text{where } f(\mathbf{v}, \mathbf{s}, \boldsymbol{\mu}) := \boldsymbol{\mu}^\top \left[(\hat{\mathbf{I}} - \mathbf{P})\mathbf{v} + \mathbf{D}\mathbf{s} \right] - \mathbf{1}^\top \mathbf{s}.$$

Note in the formulation we pose the additional constraints on \mathbf{v}, \mathbf{s} for the sake of analysis. These constraints don't change the problem optimality by noticing $\mathbf{v}^* \in \mathbb{B}_{2M}^{\mathcal{S}}, \mathbf{s}^* \in \Delta^K$; see Appendix G for details.

By designing gradient estimators and choosing divergence terms properly, one can obtain an approximately optimal solution efficiently, and thus an approximately optimal policy.

Corollary 4. *Given mixing AMDP tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathbf{r})$ with constraints $D := \max_{i,a_i,k} |[d_k]_{i,a_i}|$, for accuracy $\epsilon \in (0, 1)$, Algorithm 1 with parameter choice $\eta^v = O(\epsilon)$, $\eta^s = O(\epsilon K^{-1} D^{-2})$, $\eta^\mu = O(\epsilon t_{\text{mix}}^{-2} D^{-2} A_{\text{tot}}^{-1})$ outputs an expected ϵ -approximate solution to constrained mixing minimax problem (15) with sample complexity $O((t_{\text{mix}}^2 A_{\text{tot}} + K) D^2 \epsilon^{-2} \log(A_{\text{tot}}))$.*

Following the similar rounding technique as in Section 5, one can show the approximate solution $\boldsymbol{\mu}^\epsilon$ for minimax problem (15) leads to an approximately optimal policy π^ϵ an of problem (13).

Corollary 5. *Following the setting of Corollary 4, the policy π^ϵ induced by the unique decomposition of $\boldsymbol{\mu}^\epsilon$ from the output satisfying $\mu_{i,a_i}^\epsilon = \lambda_i^\epsilon \cdot \pi_{i,a_i}^\epsilon$, is an $O(\epsilon)$ -approximate policy for system (13).*

7 Conclusion

This work offers a general framework based on stochastic mirror descent to find an ϵ -optimal policy for AMDPs and DMDPs. It improves over previous convex optimization approaches for solving these MDPs, achieving a better sample complexity and removing an ergodicity condition for mixing AMDP, while matching the known nearly-optimal algorithms up to $(1 - \gamma)^{-1}$ factor for DMDPs.

This work reveals an interesting connection MDP problems and ℓ_∞ -regression. We believe there are a number of interesting directions and open problems for future work, including getting optimal sample complexity for discounted case, obtaining high-precision algorithms, extending the framework to broader classes of MDPs, etc. See Appendix F for a more detailed discussion of these open directions. We hope a better understanding of these problems could lead to a more complete picture of solving MDP and RL using convex-optimization methods.

Acknowledgements

This research was partially supported by NSF CAREER Award CCF-1844855, a PayPal research gift award, and a Stanford Graduate Fellowship. We thank Yair Carmon and Kevin Tian for helpful discussions on coordinate methods for matrix games; we thank Mengdi Wang, Xian Wu, Lin F. Yang, and Yinyu Ye for helpful discussions regarding MDPs; we also thank the anonymous reviewers who helped improve the completeness and readability of this paper by providing many helpful comments.

References

- Agarwal, A., Kakade, S., and Yang, L. F. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pp. 67–83, 2020.
- Altman, E. *Constrained Markov decision processes*, volume 7. CRC Press, 1999.
- Auer, P. and Ortner, R. Logarithmic online regret bounds for undiscounted reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 49–56, 2007.
- Azar, M. G., Munos, R., and Kappen, B. On the sample complexity of reinforcement learning with a generative model. *arXiv preprint arXiv:1206.6461*, 2012.
- Bertsekas, D. P. and Tsitsiklis, J. N. Neuro-dynamic programming: an overview. In *Proceedings of 1995 34th IEEE Conference on Decision and Control*, volume 1, pp. 560–564. IEEE, 1995.
- Bertsekas, D. P., Bertsekas, D. P., Bertsekas, D. P., and Bertsekas, D. P. *Dynamic programming and optimal control*, volume 1. Athena scientific Belmont, MA, 1995.
- Carmon, Y., Jin, Y., Sidford, A., and Tian, K. Variance reduction for matrix games. In *Advances in Neural Information Processing Systems*, pp. 11377–11388, 2019.
- Carmon, Y., Jin, Y., Sidford, A., and Tian, K. Coordinate methods for matrix games. *to appear in Symposium on Foundations of Computer Science*, 2020.
- Clarkson, K. L., Hazan, E., and Woodruff, D. P. Sublinear optimization for machine learning. *Journal of the ACM (JACM)*, 59(5):1–49, 2012.
- Cohen, M. B., Kelner, J., Peebles, J., Peng, R., Sidford, A., and Vladu, A. Faster algorithms for computing the stationary distribution, simulating random walks, and more. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 583–592. IEEE, 2016.
- Grigoriadis, M. D. and Khachiyan, L. G. A sublinear-time randomized approximation algorithm for matrix games. *Operations Research Letters*, 18(2):53–58, 1995.
- Kakade, S. M. et al. *On the sample complexity of reinforcement learning*. PhD thesis, University of London London, England, 2003.
- Kearns, M. J. and Singh, S. P. Finite-sample convergence rates for q-learning and indirect algorithms. In *Advances in neural information processing systems*, pp. 996–1002, 1999.
- Lee, Y. T. and Sidford, A. Path finding methods for linear programming: Solving linear programs in $o(\text{vrnk})$ iterations and faster algorithms for maximum flow. In *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*, pp. 424–433. IEEE, 2014.
- Lee, Y. T. and Sidford, A. Efficient inverse maintenance and faster algorithms for linear programming. In *2015 IEEE 56th Annual Symposium on Foundations of Computer Science*, pp. 230–249. IEEE, 2015.
- Mahadevan, S. Average reward reinforcement learning: Foundations, algorithms, and empirical results. *Machine learning*, 22(1-3):159–195, 1996.
- Nemirovski, A. Prox-method with rate of convergence $o(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251, 2004a.
- Nemirovski, A. Prox-method with rate of convergence $o(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1): 229–251, 2004b.
- Nesterov, Y. Smooth minimization of non-smooth functions. *Mathematical programming*, 103(1):127–152, 2005.
- Nesterov, Y. Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming*, 109(2-3):319–344, 2007.
- Ouyang, Y., Gagrani, M., Nayyar, A., and Jain, R. Learning unknown markov decision processes: A thompson sampling approach. In *Advances in Neural Information Processing Systems*, pp. 1333–1342, 2017.
- Palaniappan, B. and Bach, F. Stochastic variance reduction methods for saddle-point problems. In *Advances in Neural Information Processing Systems*, pp. 1416–1424, 2016.

- Puterman, M. L. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Shalev-Shwartz, S. et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.
- Sherman, J. Area-convexity, l-infinity regularization, and undirected multicommodity flow. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 452–460, 2017.
- Sidford, A. and Tian, K. Coordinate methods for accelerating l-infinity regression and faster approximate maximum flow. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 922–933. IEEE, 2018.
- Sidford, A., Wang, M., Wu, X., Yang, L., and Ye, Y. Near-optimal time and sample complexities for solving markov decision processes with a generative model. In *Advances in Neural Information Processing Systems*, pp. 5186–5196, 2018a.
- Sidford, A., Wang, M., Wu, X., and Ye, Y. Variance reduced value iteration and faster algorithms for solving markov decision processes. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 770–787. SIAM, 2018b.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Wainwright, M. J. Variance-reduced q-learning is minimax optimal. *arXiv preprint arXiv:1906.04697*, 2019.
- Wang, M. Randomized linear programming solves the discounted markov decision problem in nearly-linear running time. *arXiv preprint arXiv:1704.01869*, 2017a.
- Wang, M. Primal-dual pi learning: Sample complexity and sublinear run time for ergodic markov decision problems. *arXiv preprint arXiv:1710.06100*, 2017b.
- Ye, Y. The simplex and policy-iteration methods are strongly polynomial for the markov decision problem with a fixed discount rate. *Mathematics of Operations Research*, 36(4):593–603, 2011.