## POLICY MIRROR DESCENT FOR REGULARIZED REINFORCEMENT LEARNING: A GENERALIZED FRAMEWORK WITH LINEAR CONVERGENCE\*

WENHAO ZHAN', SHICONG CEN', BAIHE HUANG<sup>S</sup>, YUXIN CHEN<sup>P</sup>,

JASON D. LEE', AND YUEJIE CHI'

Abstract. Policy optimization, which learns the policy of interest by maximizing the value function via large-scale optimization techniques, lies at the heart of modern reinforcement learning (RL). In addition to value maximization, other practical considerations arise commonly as well, including the need of encouraging exploration, and that of ensuring certain structural properties of the learned policy due to safety, resource, and operational constraints. These considerations can often be accounted for by resorting to regularized RL, which augments the target value function with a structure-promoting regularization term. Focusing on an infinite-horizon discounted tabular Markov decision process, this paper proposes a generalized policy mirror descent (GPMD) algorithm for solving regularized RL. As a generalization of policy mirror descent [G. Lan, Math. Program., 198 (2023), pp. 1059--1106], the proposed algorithm accommodates a general class of convex regularizers as well as a broad family of Bregman divergence in cognizance of the regularizer in use. We demon-strate that our algorithm converges linearly to the global solution over an entire range of learning rates, in a dimension-free fashion, even when the regularizer lacks strong convexity and smoothness. In addition, this linear convergence feature is provably stable in the face of inexact policy evaluation and imperfect policy updates. Numerical experiments are provided to corroborate the applicability and appealing performance of GPMD.

Key words. policy mirror descent, Bregman divergence, regularization, policy optimization

MSC codes. 68Q25, 68Q32, 90C26

DOI. 10.1137/21M1456789

1. Introduction. Policy optimization lies at the heart of recent successes of reinforcement learning (RL) [39]. In its basic form, the optimal policy of interest, or a suitably parameterized version, is learned by attempting to maximize the value

<sup>\*</sup>Received by the editors November 2, 2021; accepted for publication (in revised form) January 6, 2023; published electronically June 22, 2023. A preliminary version of this work was presented at NeurIPS 2021 Workshop on Optimization for Machine Learning. The first two authors contributed equally.

https://doi.org/10.1137/21M1456789

Funding: The research of the second and sixth authors was supported in part by grants ONR N00014-19-1-2404, NSF CCF-2106778, DMS-2134080, CCF-1901199, CCF-2007911, and CNS-2148212. The research of the second author was also supported by Wei Shen and Xuehong Zhang's Presidential Fellowship and Nicholas Minnici Dean's Graduate Fellowship in Electrical and Computer Engineering at Carnegie Mellon University. The research of the first and fourth authors was supported in part by the Google Research Scholar Award, the Alfred P. Sloan Research Fellowship, and grants AFOSR FA9550-22-1-0198, ONR N00014-22-1-2354, NSF CCF-2221009, CCF-1907661, IIS-2218713, and IIS-2218773. The research of the first and fifth authors was supported in part by the ARO under MURI award W911NF-11-1-0304, the Sloan Research Fellowship, NSF CCF 2002272, NSF IIS 2107304, and an ONR Young Investigator Award.

Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ 08544 USA (wz3993@princeton.edu, jasondl@princeton.edu).

Department of Electrical and Computer Engineering, Carnegie Mellon University, Pittsburgh, PA 15213 USA (shicongc@andrew.cmu.edu, yuejiechi@cmu.edu).

<sup>&</sup>lt;sup>S</sup> Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, Berkeley, CA 94720 USA (baihehuang@pku.edu.cn).

PDepartment of Statistics and Data Science, Wharton School, University of Pennsylvania, Philadelphia, PA 19104 USA (yuxinc@wharton.upenn.edu).

function in a Markov decision process (MDP). For the most part, the maximization step is carried out by means of first-order optimization algorithms amenable to large-scale applications, whose foundations were set forth in the early works of [53, 47]. A partial list of widely adopted variants in modern practice includes policy gradient (PG) methods [47], natural policy gradient (NPG) methods [24], TRPO [44], PPO [45], and soft actor-critic methods [22], to name just a few. In comparison with model-based and value-based approaches, this family of policy-based algorithms offers a remarkably flexible framework that accommodates both continuous and discrete action spaces and lends itself well to the incorporation of powerful function approximation schemes like neural networks. In stark contrast to its practical success, however, theoretical understanding of policy optimization remains severely limited even for the tabular case, largely owing to the ubiquitous nonconvexity issue underlying the objective function.

- 1.1. The role of regularization. In practice, there are often competing objectives and additional constraints that the agent has to deal with in conjunction with maximizing values, which motivate the studies of regularization techniques in RL. In what follows, we isolate a few representative examples.
  - t Promoting exploration. In the face of large problem dimensions and complex dynamics, it is often desirable to maintain a suitable degree of randomness in the policy iterates, in order to encourage exploration and discourage premature convergence to suboptimal policies. A popular strategy of this kind is to enforce entropy regularization [54], which penalizes policies that are not suficiently stochastic. Along similar lines, the Tsallis entropy regularization [16, 31] further promotes sparsity of the learned policy while encouraging exploration, ensuring that the resulting policy does not assign non-negligible probabilities to too many suboptimal actions.
  - t Safe RL. In a variety of application scenarios such as industrial robot arms and self-driving vehicles, the agents are required to operate safely both to themselves and the surroundings [5, 40]; for example, certain actions might be strictly forbidden in some states. One way to incorporate such prescribed operational constraints is through adding a regularizer (e.g., a properly chosen log barrier or indicator function tailored to the constraints) to explicitly account for the constraints.
  - t Cost-sensitive RL. In reality, different actions of an agent might incur drastically different costs even for the same state. This motivates the design of new objective functions that properly trade off the cumulative rewards against the accumulated cost, which often take the form of certain regularized value functions.

Viewed in this light, it is of imminent value to develop a unified framework towards understanding the capability and limitations of regularized policy optimization. While a recent line of works [3, 38, 12] has looked into specific types of regularization techniques such as entropy regularization, existing convergence theory remains highly inadequate when it comes to a more general family of regularizers.

1.2. Main contributions. The current paper focuses on policy optimization for regularized RL in a a-discounted infinite horizon Markov decision process (MDP) with state space S, action space A, and reward function r(t,t). The goal is to find an optimal policy that maximizes a regularized value function. Informally speaking, the regularized value function associated with a given policy i takes the following form:

$$V_u^i = V^i - uE \begin{bmatrix} (t | s) \end{bmatrix}$$

where  $V^i$  denotes the original (unregularized) value function, u>0 is the regularization parameter,  $h_s(t)$  denotes a convex regularizer employed to regularize the policy in state s, and the expectation is taken over certain marginal state distribution w.r.t. the MDP (to be made precise in section 2.1). It is noteworthy that this paper does not require the regularizer  $h_s$  to be either strongly convex or smooth.

In order to maximize the regularized value function (2.8b), Lan [27] exhibited a seminal algorithm called Policy Mirror Descent (PMD), which can be viewed as an adaptation of the mirror descent algorithm [41, 7] to the realm of policy optimization. In particular, PMD subsumes the natural policy gradient (NPG) method [24] as a special case. To further generalize PMD [27], we propose an algorithm called Generalized Policy Mirror Descent (GPMD). In each iteration, the policy is updated for each state in parallel via a mirror-descent style update rule. In sharp contrast to [27], which considered a generic Bregman divergence, our algorithm selects the Bregman divergence adaptively in cognizance of the regularizer, which leads to complementary perspectives and insights. Several important features and theoretical appeal of GPMD are summarized as follows.

- t GPMD substantially broadens the range of (provably effective) algorithmic choices for regularized RL and subsumes several well-known algorithms as special cases. For example, it reduces to regularized policy iteration [21] when the learning rate tends to infinity and subsumes entropy-regularized NPG methods as special cases if we take the Bregman divergence to be the Kullback--Leibler (KL) divergence [12].
- t Assuming exact policy evaluation and perfect policy update in each iteration, GPMD converges linearly--in a dimension-free fashion--over the entire range of the learning rate a > 0. More precisely, it converges to an n-optimal regularized Q-function in no more than an order of

$$\frac{1+au}{au(1-a)}\log\frac{1}{n}$$

iterations (up to some logarithmic factor). Encouragingly, this appealing feature is valid for a broad family of convex and possibly nonsmooth regularizers.

- t The intriguing convergence guarantees are robust in the face of inexact policy evaluation and imperfect policy updates; namely, the algorithm is guaranteed to converge linearly at the same rate until an error floor is hit. See section 3.2 for details.
- t Numerical experiments are provided in section 5 to demonstrate the practical applicability and appealing performance of the proposed GPMD algorithm.

Finally, we find it helpful to briefly compare the above findings with prior works. As soon as the learning rate exceeds a q 1/u, the iteration complexity of our algo-rithm is at most on the order of  $\frac{1}{1}$  log  $\frac{1}{1}$ , thus matching that of regularized policy iteration [21]. In comparison to [27], our work sets forth a different framework to an-alyze mirror-descent-type algorithms for regularized policy optimization, generalizing and refining the approach in [12] far beyond entropy regularization. When constant learning rates are employed, the linear convergence of PMD [27] critically requires the regularizer to be strongly convex, with only sublinear convergence theory established for convex regularizers. In contrast, we establish the linear convergence of GPMD under constant learning rates even in the absence of strong convexity. Furthermore, for the special case of entropy regularization, the stability analysis of GPMD also significantly improves over the prior art in [12], preventing the error floor from blowing up when the learning rate approaches zero, as well as incorporating the impact of

optimization error that was previously uncaptured. More detailed comparisons with [27, 12] can be found in section 3.

1.3. Related works. Before embarking on our algorithmic and theoretic developments, we briefly review a small sample of other related works.

Global convergence of policy gradient methods. Recent years have witnessed a surge of activities towards understanding the global convergence properties of policy gradient methods and their variants for both continuous and discrete RL problems, examples including [20, 9, 3, 63, 51, 37, 10, 25, 35, 37, 4, 56, 51, 12, 36, 34, 52, 60, 62, 61, 46], among other things. The authors of [42] provided the first interpretation of NPG methods as mirror descent [41], thereby enabling the adaptation of techniques for analyzing mirror descent to the studies of NPG-type algorithms such as TRPO [46, 48]. It has been shown that the NPG method converges sublinearly for unregularized MDPs with a fixed learning rate [3], and converges linearly if the learning rate is set adaptively [25], via exact line search [10], or following a geometrically increasing schedule [55]. The global linear convergence of NPG holds more generally for an arbitrary fixed learning rate when entropy regularization is enforced [12]. Noteworthily, the authors of [33] established a lower bound indicating that softmax PG methods can take an exponential time--in the size of the state space--to converge, while the convergence rates of NPG-type methods are almost independent of the problem dimension. In addition, another line of recent works [1, 23, 30] established regret bounds for approximate NPG methods--termed as KL-regularized approximate policy iteration therein--for infinite-horizen undiscounted MDPs, which are beyond the scope of the current paper.

Regularization in RL. Regularization has been suggested to the RL literature either through the lens of optimization [17, 3] or through the lens of dynamic programming [21, 50]. Our work is clearly an instance of the former type. Several recent results in the literature merit particular attention: The authors of [3] demonstrated sublinear convergence guarantees for PG methods in the presence of relative entropy regularization and the authors of [38] established linear convergence of entropy-regularized PG methods, whereas the authors of [12] derived an almost dimension-free linear convergence theory for NPG methods with entropy regularization. Most of the existing literature focused on the entropy regularization or KL-type regularization, and the studies of general regularizers had been quite limited until the recent work [27]. The regularized MDP problems are also closely related to the studies of constrained MDPs, as both types of problems can be employed to model/promote constraint satisfaction in RL, as recently investigated in, e.g., [15, 19, 18, 58, 57]. Note, however, that it is dificult to directly compare our algorithm with these methods, due to drastically different formulations and settings.

1.4. Notation. Let us introduce several notations that will be adopted throughout. For any set A, we denote by |A| the cardinality of a set A and let a (A) indicate the probability simplex over the set A. For any convex and differentiable function h(t), the Bregman divergence generated by h(t) is defined as

(1.1) 
$$D_h(z,x) := h(z) - h(x) - a h(x), z - x.$$

For any convex (but not necessarily differentiable) function h(t), we denote by I h the subdifferential of h. Given two probability distributions i  $_1$  and i  $_2$  over A , the K L divergence from i  $_2$  to i  $_1$  is defined as K L(i  $_1$  | i  $_2$ ) :=  $_{an\ A}^{n}$  i  $_1$ (a) log  $_{i\ 2}^{i\ 1}$ (a). For any vectors a =  $[a_i]_{1q\ iq\ n}$  and b =  $[b_i]_{1q\ iq\ n}$ , the notation a q b (resp., a q b) means that

 $a_i \neq b_i$  ( $a_i \neq b_i$ ) for every 1  $\neq i \neq n$ . We shall also use 1 (resp., 0) to denote the all-one (resp., all-zero) vector whenever it is clear from the context.

## Model and algorithms.

For any policy i, we define the associated value function  $V^{i}:S$  w R as follows:

(2.1) 
$$| s n s : V^{i}(s) := E \begin{cases} W & | W \\ atmi(t|st), & | t=0 \end{cases}$$

$$t = 0$$

$$t = 0$$

which can be viewed as the utility function we wish to maximize. Here, the expectation is taken over the randomness of the MDP trajectory  $\{(s_t, a_t)\}_{tq\ 0}$  induced by policy i . Similarly, when the initial action a is fixed, we can define the action-value function (or Q-function) as follows:

A well-known fact is that the policy gradient of  $V^i$  (w.r.t. the policy i ) admits the following closed-form expression [47]:

(2.3) 
$$I(s,a) \cap S S A : \frac{IV^{i}(s_{0})}{Ii(a|s)} = \frac{1}{1-a} d_{s_{0}}^{i}(s) Q^{i}(s,a).$$

Here,  $d_{s_0}^i$  n a (S ) is the so-called discounted state visitation distribution defined as follows:

(2.4) 
$$d_{s_0}^i(s) := (1 - a) \prod_{t=0}^{n} a^t P^i(s_t = s | s_0),$$

where  $P^i$  ( $s_t = s | s_0$ ) denotes the probability of  $s_t = s$  when the MDP trajectory  $\{s_t\}_{t \neq 0}$  is generated under policy i given the initial state  $s_0$ .

Furthermore, the optimal value function and the optimal Q-function are defined and denoted by

(2.5) I (s, a) n S s A : 
$$V^{r}(s) := \max_{i} V^{i}(s), Q^{r}(s, a) := \max_{i} Q^{i}(s, a).$$

It is well known that there exists at least one optimal policy, denoted by i<sup>r</sup>, that simultaneously maximizes the value function and the Q-function for all state-action pairs [2].

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

Regularized MDP. In practice, the agent is often asked to design policies that possess certain structural properties in order to be cognizant of system constraints such as safety and operational constraints, as well as to encourage exploration during the optimization/learning stage. A natural strategy to achieve these is to resort to the following regularized value function w.r.t. a given policy i [42, 38, 12, 27]:

where h<sub>s</sub>: a <sub>a</sub>(A) w R stands for a convex and possibly nonsmooth regularizer for state s, u > 0 denotes the regularization parameter, and  $d_s^i(t)$  is defined in (2.4). Here, for technical convenience, we assume throughout that  $h_s(t)$  (s n S) is well-defined over an "a -neighborhood"" of the probability simplex a (A ) defined as follows:

where a > 0 can be an arbitrary constant. For instance, entropy regularization adopts the choice  $h_s(p) = \prod_{i \in A}^n p_i \log p_i$  for all s n S and p n a (A ), which coincides with the negative Shannon entropy of a probability distribution. Similarly, a K L regularization adopts the choice  $h_s(p) = KL(p \mid p_{ref})$ , which penalizes the distribution p that deviates from the reference  $p_{ref}$ . As another example, a weighted  $I_1$  regularization adopts the  $w_{s,i}p_i$  for all s n S and p n a (A ), where  $w_{s,i} \neq 0$  is the cost of taking action i at state s, and the regularizer  $h_s(i \mid t \mid s)$  captures the expected cost of the policy i in state s. Throughout this paper, we impose the following assumption.

Assumption 1. Consider an arbitrarily small constant a > 0. For for any s n S, suppose that h<sub>s</sub>(t) is convex and

(2.7) 
$$h_s(p) = y \qquad \text{for any } p \eta' a_a(A).$$

Following the convention in prior literature (see, e.g., [38]), we also define the corresponding regularized Q-function as follows:

(2.8a) I (s, a) n S s A : 
$$Q_{\iota}^{i}(s, a) := r(s, a) + a E_{s^{e} \text{ mP } (t \mid s, a)} [V_{\iota}^{i}(s^{e})].$$

As can be straightforwardly verified, one can also express 
$$V_u^i$$
 in terms of  $Q_u^i$  as (2.8b) IsnS:  $V_u^i$ (s):=  $E_{am\ i\ (t\ |s)}^i Q_u^i$ (s,a)-  $uh_s^i$ (t|s).

The optimal regularized value function  $V_u^r$  and the corresponding optimal policy i  $_u^r$ are defined, respectively, as follows:

(2.9) IsnS: 
$$V_u^r(s) := V_u^{i_u^r}(s) = \max_i V_u^i(s), \quad i_u^r := \arg\max_i V_u^i.$$

It is worth noting that the author of [43] asserts the existence of an optimal policy i that achieves (2.9) simultaneously for all s n S. Correspondingly, we shall also define the resulting optimal regularized Q-function as

2.2. Algorithm: Generalized policy mirror descent. Motivated by PMD [27], we put forward a generalization of PMD that selects the Bregman divergence in cognizance of the regularizer in use. A thorough comparison with [27] will be provided after introducing our generalized PMD algorithm.

For notational simplicity, we shall write

$$(2.11) \quad V_u^{(k)} := V_u^{i^{(k)}}, \qquad Q_u^{(k)}(s,a) := Q_u^{i^{(k)}}(s,a), \qquad \text{and} \qquad d_{s_0}^{(k)}(s) := d_{s_0}^{i^{(k)}}(s)$$

throughout the paper, where i (k) denotes our policy estimate in the kth iteration.

To begin with, suppose for simplicity that  $h_s(t)$  is differentiable everywhere. In the kth iteration, a natural MD scheme that comes into mind for solving (2.6)-namely, maximize;  $V_u^i$  (s<sub>0</sub>) for a given initial state s<sub>0</sub>--is the following update rule:

(2.12)

$$\begin{split} &i^{(k+1)}(t\mid s) \\ &= \underset{pn \, a \, (A)}{\text{min}} \left\{ \begin{array}{l} -e \\ -e \\ a_{i \, (t\mid s)} V_u^{i} \left(s_0\right) \Big|_{i \, = \, i^{\, (k)}}, p \end{array} \right. \\ &+ \frac{l}{1 - a} d_{s_0}^{(k)}(s) h_s(p) + \frac{1}{a^e} D_{h_s} \left(p, i^{\, (k)}(t\mid s)\right)^{s} \\ &= \underset{pn \, a \, (A)}{\text{min}} \left\{ \begin{array}{l} \frac{1}{1 - a} d_{s_0}^{(k)}(s)^{s} - e_{Q_t^{(k)}(s,t), p} + uh_s(p) + \frac{1}{a^e} D_{h_s} \left(p, i^{\, (k)}(t\mid s)\right)^{s} \\ &= \underset{pn \, a \, (A)}{\text{min}} - e_{Q_t^{(k)}(s,t), p} + uh_s(p) + \frac{1}{\epsilon} D_{h_s} \left(p, i^{\, (k)}(t\mid s)\right)^{s} \end{array} \right. \end{split}$$

for every state s n S . Here, we start with a learning rate  $a^e$  and obtain simplification by replacing  $a^e$  with  $a(1 - a)/d_{s_0}^{(k)}(s)$ . Notably, the update strategy (2.12) is invariant to the initial state  $s_0$ , akin to natural policy gradient methods [3].

This update rule is well-defined for, say, the case when  $h_s$  is the negative entropy, since the algorithm guarantees i <sup>(k)</sup> > 0 all the time and hence  $h_s$  is always differentiable w.r.t. the kth iterate (see [12]). In general, however, it is possible to encounter situations when the gradient of  $h_s$  does not exist on the boundary (e.g., when  $h_s$  represents a certain indicator function). To cope with such cases, we resort to a generalized version of Bregman divergence (see, e.g., [26, 28, 29]). To be specific, we attempt to replace the usual Bregman divergence  $D_{h_s}$  (p, q) by the following metric:

(2.13) 
$$D_{h_s}(p,q;g_s) := h_s(p) - h_s(q) - eg_s, p - qe q 0,$$

where  $g_s$  can be any vector falling within the subdifferential I  $h_s(q)$ . Here, the non-negativity condition in (2.13) follows directly from the definition of the subgradient for any convex function. The constraint on  $g_s$  can be further relaxed by exploiting the requirement p,q n a (A). In fact, for any vector  $i_s = g_s - c_s 1$  (with  $c_s$  n R some constant and 1 the all-one vector), one can readily see that

$$D_{h_s}(p,q;g_s) = h_s(p) - h_s(q) - eg_s, p - qe = h_s(p) - h_s(q) - ei_s, p - qe + c_se 1, p - qe$$
  
(2.14) =  $h_s(p) - h_s(q) - ei_s, p - qe = D_{h_s}(p,q;i_s),$ 

where the last line is valid since  $1^p$  p =  $1^p$  q = 1. As a result, everything boils down to identifying a vector  $i_s$  that falls within  $I h_s(q)$  upon global shift.

Towards this, we propose the following iterative rule for designing such a sequence of vectors as surrogates for the subgradient of h<sub>s</sub>:

(2.15a) 
$$i^{(0)}(s,t) n I h_s^{(i^{(0)}(t \mid s))};$$

(2.15b) 
$$i^{(k+1)}(s,t) = \frac{1}{1+au}i^{(k)}(s,t) + \frac{a}{1+au}Q_t^{(k)}(s,t), \quad k \neq 0,$$

where i  $^{(k+1)}(s,t)$  is updated as a convex combination of the previous i  $^{(k)}(s,t)$  and  $Q_t^{(k)}(s,t)$ , where more emphasis is put on  $Q_t^{(k)}(s,t)$  when the learning rate a is large. As asserted by the following lemma, the above vectors i  $^{(k)}(s,t)$  we construct satisfy the desired property, i.e., lying within the subdifferential of  $h_s$  under suitable global shifts. It is worth mentioning that these global shifts  $\{c_s^{(k)}\}$  only serve as an aid to better understand the construction but are not required during the algorithm updates.

Lemma 2.1. For all k q 0 and every s n S, there exists a quantity  $c_s^{(k)} \, n \, R$  such that

(2.16) 
$$i^{(k)}(s,t) - c_s^{(k)} 1 \text{ n I h}_s^{(k)} i^{(k)}(t \mid s)^{(k)}$$

In addition, for every s n S , there exists a quantity  $c_s^{r} \; n \; R \; \text{such that}$ 

(2.17) 
$$u^{-1}Q_{t}^{r}(s,t) - c_{s}^{r} 1 n | h_{s}^{r} i_{u}^{r}(t|s)$$

Proof. See Appendix A.1.

Thus far, we have presented all crucial ingredients of our algorithm. The whole procedure is summarized in Algorithm 2.1 and will be referred to as Generalized Policy Mirror Descent (GPMD) throughout the paper. Interestingly, several well-known algorithms can be recovered as special cases of GPMD.

## Algorithm 2.1. PMD with generalized Bregman divergence (GPMD).

Input: initial policy iterate i  $^{(0)}$ , learning rate a > 0.

Initialize i  $^{(0)}$  so that i  $^{(0)}(s,t)$  n I  $h_s(i^{(0)}(t\mid s))$  for all s n S .

for k = 0, 1, ..., do

For every s n S, set

$$i^{(k+1)}(t|s) = \arg\min_{p \in a \text{ (A)}} \left\{ e_{t}^{(k)}(s,t), p + uh_{s}(p) + \frac{1}{\epsilon}D_{h_{s}}^{(k)}(t|s); i^{(k)} \right\},$$

where

(2.18b) 
$$D_{h_s} (p,q;i) := h_s(p) - h_s(q) - i(s,t), p-q.$$

For every (s, a) n S s A, compute

(2.18c) 
$$i^{(k+1)}(s,a) = \frac{1}{1+au}i^{(k)}(s,a) + \frac{\tilde{\epsilon}}{1+au}Q_{\iota}^{(k)}(s,a).$$

end for

- t When the Bregman divergence  $D_{h_s}(t,t)$  is taken as the K L divergence, GPMD reduces to the well-renowned NPG algorithm [24] when u=0 (no regularization) and to the NPG algorithm with entropy regularization analyzed in [12] when  $h_s(t)$  is taken as the negative Shannon entropy.
- t When a = y (no divergence), GPMD reduces to regularized policy iteration in [21]; in particular, GPMD reduces to the standard policy iteration algorithm if in addition u is also 0.

Comparison with PMD [27]. Before continuing, let us take a moment to point out the key differences between our algorithm GPMD and the PMD algorithm proposed in [27] in terms of algorithm designs. Although the primary exposition of PMD in [27] fixes the Bregman divergence as the KL divergence, the algorithm also works in the presence of a generic Bregman divergence, whose relationship with the regularizer  $h_s$  is, however, unspecified. Furthermore, GPMD adaptively sets this term to be the Bregman divergence generated by the regularizer  $h_s$  in use, together with a carefully designed recursive update rule (cf. (2.15)) to compute surrogates for the subgradient of  $h_s$  to facilitate implementation. Encouragingly, this specific choice leads to a tailored performance analysis of GPMD, which was not present in and instead complementary with that of PMD [27]. In truth, our theory offers linear convergence guarantees for more general scenarios by adapting to the geometry of the regularizer  $h_s$ ; details will follow momentarily.

- 3. Main results. This section presents our convergence guarantees for the GPMD method presented in Algorithm 2.1. We shall start with the idealized case, assuming that the update rule can be precisely implemented, and then discuss how to generalize it to the scenario with imperfect policy evaluation.
- 3.1. Convergence of exact GPMD. To start with, let us pin down the convergence behavior of GPMD, assuming that accurate evaluation of the policy  $Q^{(k)}$  is available and the subproblem (2.18a) can be solved perfectly. Here and below, we shall refer to the algorithm in this case as exact GPMD. Encouragingly, exact GPMD provably achieves global linear convergence from an arbitrary initialization, as asserted by the following theorem.

Theorem 3.1 (exact GPMD). Suppose that Assumption 1 holds. Consider any learning rate a > 0, and set a :=  $\frac{1}{1+a}$ . Then the iterates of Algorithm 2.1 satisfy

for all k q 0, where  $C_1:=\,|\,Q_u^r\,-\,\,Q_t^{(0)}\,|_{\,y}\,\,+\,2a\,\,|\,Q_u^r\,-\,\,ui\,^{(0)}\,|_{\,y}\,$  .

In addition, if  $h_s$  is 1-strongly convex w.r.t. the  $l_1$  norm for some s n S, then one further has

(3.2) 
$$\left| i_{u}^{r}(s) - i_{u}^{(k+1)}(s) \right|_{1}^{l} q u^{-1} (1 - a)(1 - a)^{k} C_{1}, \quad kq 0.$$

Our theorem confirms the fast global convergence of the GPMD algorithm, in terms of both the resulting regularized Q-value (if  $h_s(t)$  is convex) and the policy estimate (if  $h_s(t)$  is strongly convex). In summary, it takes GPMD no more than

(3.3a) 
$$\frac{1}{(1-a)(1-a)}\log\frac{C_1}{n} = \frac{1+au}{au(1-a)}\log\frac{C_1}{n}$$

iterations to converge to an n-optimal regularized Q-function (in the I, sense) or

(3.3b) 
$$\frac{1}{(1-a)(1-a)}\log\frac{C_1}{nu} = \frac{1+au}{au(1-a)}\log\frac{C_1}{nu}$$

iterations to yield an n-approximation (w.r.t. the  $I_1$  norm error) of  $i_1^r$ . The iteration complexity (3.3) is nearly dimension-free--namely depending at most logarithmically on the dimension of the state-action space--making it scalable to large-dimensional problems.

To make clear our contributions, it is helpful to compare Theorem 3.1 with the theory for the state-of-the-art algorithm PMD in [27].

- t Linear convergence for convex regularizers under constant learning rates. Suppose that constant learning rates are adopted for both GPMD and PMD. Our finding reveals that GPMD enjoys global linear convergence--in terms of both  $|\,Q_t^r\,-\,Q_t^{(k+1)}\,|_y\,$  and  $|\,V_t^r\,-\,V_u^{(k+1)}\,|_y\,$  --even when the regularizer  $h_s(t)$  is only convex but not strongly convex. In contrast, [27, Theorem 2] provided only sublinear convergence guarantees (with an iteration complexity proportional to 1/n) for the case with convex regularizers, provided that constant learning rates are adopted.  $^1$
- t A full range of learning rates. Theorem 3.1 reveals linear convergence of GPMD for a full range of learning rates; namely, our result is applicable to any a > 0. In comparison, linear convergence was established in [27] only when the learning rates were suficiently large and when h<sub>s</sub> was 1-strongly convex w.r.t. the KL divergence. Consequently, the linear convergence results in [27] do not extend to several widely used regularizers, such as negative Tsallis entropy and log-barrier functions (even after scaling), which are, in contrast, covered by our theory. It is worth noting that the case with small-to-medium learning rates is often more challenging to cope with in theory, given that its dynamics could differ drastically from that of regularized policy iteration.
- t Further comparison of rates under large learning rates. [27, Theorem 1] achieves a contraction rate of a when the regularizer is strongly convex and the step size satisfies a q  $\frac{1-a}{a\,t}$ , while the contraction rate of GPMD is  $1-\frac{a\,u}{1+a\,t}$  (1 a) under the full range of the step size, which is slower but approaches the contraction rate a of PMD as a goes to infinity. Therefore, in the limit a w y , both GPMD and PMD achieve the contraction rate a. As soon as a q 1/u, their iteration complexities are on the same order.
- 3.2. Convergence of approximate GPMD. In reality, however, it is often the case that GPMD cannot be implemented in an exact manner, either because perfect policy evaluation is unavailable or because the subproblem (2.18a) cannot be solved exactly. To accommodate these practical considerations, this subsection generalizes our previous result by permitting inexact policy evaluation and non-zero optimization error in solving (2.18a). The following assumptions make precise this imperfect scenario.

Assumption 2 (policy evaluation error). Suppose, for any k q 0, we have access to an estimate  $\mathbb{Q}_{+}^{(k)}$  obeying

$$\left| Q_{u}^{(k)} - Q_{t}^{(k)} \right|_{y} \neq n_{\text{eval}}.$$

<sup>&</sup>lt;sup>1</sup>In fact, [27, Theorem 3] suggests using a vanishing strongly convex regularization, as well as a corresponding increasing sequence of learning rates, in order to enable linear convergence for non-strongly-convex regularizers.

Algorithm 3.1. Approximate PMD with generalized Bregman divergence (approximate GPMD).

Input: initial policy i (0), learning rate a > 0.

Initialize  $t^{(0)}(s) \cap I h_s(i^{(0)}(t \mid s))$  for all  $s \cap S$ .

for k = 0, 1, ..., do

For every s n S, invoke the oracle to obtain (cf. (3.5))

(3.6) 
$$i^{(k+1)}(s) = G_{s,n_{opt}} (Q_{t}^{(k)}, i^{(k)}, i^{t(k)})$$

For every (s, a) n S s A, compute

(3.7) 
$$t^{(k+1)}(s,a) = \frac{1}{1+au} t^{(k)}(s,a) + \frac{\epsilon}{1+au} Q_{\iota}^{(k)}(s,a).$$

end for

Assumption 3 (subproblem optimization error). Consider any policy i and any vector i n R $_{s}^{|}$  $_{u}^{|}$ . Define

$$f_s(p;i,i) := -\frac{e}{Q(s,t),p} + uh_s(p) + \frac{1}{a}D_{h_s}(p,i(t|s);i(s,t)),$$

where  $D_{h_s}(p,q;i)$  is defined in (2.13). Suppose there exists an oracle  $G_{s,n_{opt}}(Q,i,i)$  that is capable of returning i e(t|s) such that

(3.5) 
$$f_{s}^{(i)}(t|s); i, i) q \min_{pn \ a \ (A)} f_{s}^{(p;i,i)} + n_{opt}^{(p)}.$$

Note that the oracle in Assumption 3 can be implemented eficiently in practice via various first-order methods [6]. Under Assumptions 2 and 3, we can modify Algorithm 2.1 by replacing  $\{Q_u^{(k)}\}$  with the estimate  $\{Q_u^{(k)}\}$  and invoking the oracle  $G_{s,n}$  of (Q,i,i) to solve the subproblem (2.18a) approximately. The whole procedure, which we shall refer to as approximate GPMD, is summarized in Algorithm 3.1.

The following theorem uncovers that approximate GPMD converges linearly--at the same rate as exact GPMD--before an error floor is hit.

Theorem 3.2 (approximate GPMD). Suppose that Assumptions 1, 2, and 3 hold. Consider any learning rate a > 0. Then the iterates of Algorithm 3.1 satisfy

where a :=  $\frac{1}{1+a}$ - $\frac{1}{u}$ , C<sub>1</sub> is defined in Theorem 3.1, and

$$C_2 := \frac{1}{1 - a} \begin{bmatrix} ( & & & \\ & 2 + \frac{2a}{(1 - a)(1 - a)} & n_{eval} + & 1 + \frac{2a}{(1 - a)(1 - a)} & n_{opt} \end{bmatrix}.$$

In addition, if  $h_{\text{\tiny S}}$  is 1-strongly convex w.r.t. the  $I_{\,1}$  norm for any s n S , then we can further obtain

(3.9a) 
$$|Q_u^r - Q_u^{(k+1)}|_y$$
 q a 1- (1- a)(1- a) $^k C_1 + C_3$ ,  
(3.9b)  $|V_u^r - V_u^{(k+1)}|_y$  q (a + 2) 1- (1- a)(1- a) $^k C_1 + C_3 + (1- a)n_{opt}$ ,

(3.9c) 
$$\left| i \left[ (t \mid s) - i \right]^{(k+1)} (t \mid s) \right|_{1}^{1} q u^{-1} \left[ (1 - a)(1 - a)^{k} C_{1} + C_{3} + \frac{2a r_{opt}}{1 + au} \right]$$

where 
$$(3.10) \qquad C_3 := \frac{1}{1-a} \quad 2 + \frac{n_{ev\bar{\epsilon}} |a|}{u(1-a)} \quad n_{eval} + \quad 1 + \frac{4a}{(1-a)(1-a)} \quad n_{opt} \ .$$

In the special case where  $n_{opt} = 0$  and a = y, Algorithm 3.1 reduces to regularized policy iteration, and the convergence result can be simplified as follows:

$$\left| Q_u^r - Q_t^{(k)} \right|_y = q \epsilon^{k} \left| Q_u^r - Q_t^{(0)} \right|_y + \frac{2a n_{eval}}{(1-a)^2}.$$

In particular, when  $h_s$  is taken as the negative entropy, our result strengthens the prior result established in [12] for the approximate entropy-regularized NPG method with  $n_{opt}=0$  over a wide range of learning rates. Specifically, the error bound in [12] reads as a t  $\frac{n_{eval}}{1-a}-2+\frac{2a}{au}$ , where the second term in the brackets scales inversely with respect to a and therefore grows unboundedly as a approaches 0. In contrast, (3.9) and (3.10) suggest a bound a t  $\frac{n_{eval}}{1-a}-2+\frac{n_{eval}}{u(1-a)}$ , which is independent of the learning rate a in use and thus prevents the error bound from blowing up when the learning rate approaches 0. Indeed, our result improves over the prior art [12] whenever a q  $\frac{2(1-a)}{n_{eval}}$ .

Remark 1 (sample complexities). One might naturally ask how many samples are suficient to learn an n-optimal regularized Q-function by leveraging sample-based policy evaluation algorithms in GPMD. Notice that it is straightforward to consider an expected version of Assumption 2 as follows:

where the expectation is with respect to the randomness in policy evaluation; then the convergence results in Theorem 3.2 apply to E [|  $Q_u^r - Q_t^{(k+1)}|_{\gamma}$ ] and E [|  $i_u^r(t|s) - i_t^{(k+1)}(t|s)|_1$ ] instead. This randomized version makes it immediately amenable to combine with, e.g., the rollout-based policy evaluators in [27, section 5.1], to obtain (possibly crude) bounds on the sample complexity. We omit these straightforward developments.

Roughly speaking, approximate GPMD is guaranteed to converge linearly to an error bound that scales linearly in both the policy evaluation error  $n_{\text{eval}}$  and the optimization error  $n_{\text{opt}}$ , thus confirming the stability of our algorithm vis- $\frac{1}{2}$ -vis imperfect implementation of the algorithm. As before, our theory improves upon prior works by demonstrating linear convergence for a full range of learning rates even in the absence of strong convexity and smoothness.

4. Analysis for exact G P M D (Theorem 3.1). In this section, we present the analysis for our main result in Theorem 3.1, which follows a framework different from

[27]. Here and throughout, we shall often employ the following shorthand notation when it is clear from the context:

(4.1) 
$$i^{(k)}(s) := i^{(k)}(t \mid s) \text{ n a } (A), \qquad Q^{i}(s) := Q^{i}(s,t) \text{ n R }_{A}^{i}, \\ i^{(k)}(s) := i^{(k)}(s,t) \text{ n R }_{A}^{i}, \qquad Q^{i}_{u}(s) := Q^{i}_{u}(s,t) \text{ n R }_{A}^{i},$$

in addition to those already defined in (2.11).

4.1. Preparation: Basic facts. In this subsection, we single out a few basic results that underlie the proof of our main theorems.

Performance improvement. To begin with, we demonstrate that GPMD enjoys a sort of monotonic improvement concerning the updates of both the value function and the Q-function, as stated in the following lemma. This lemma can be viewed as a generalization of the well-established policy improvement lemma in the analysis of NPG [3, 12] as well as PMD [27]. The proof can be found in [59].

Lemma 4.1 (pointwise monotonicity). For any (s,a) n S s A and any k q 0, Algorithm 2.1 achieves

(4.2) 
$$V_u^{(k+1)}(s) \neq V_u^{(k)}(s)$$
 and  $Q_u^{(k+1)}(s,a) \neq Q_u^{(k)}(s,a)$ .

Interestingly, the above monotonicity holds simultaneously for all state-action pairs and hence can be understood as a kind of pointwise monotonicity.

Generalized Bellman operator. Another key ingredient of our proof lies in the use of a generalized Bellman operator  $T_{u,h}:R^{|S|}|_A|wR_S^{|S|}|_A$  associated with the regularizer  $h=\{h_s\}_{sn|S}$ . Specifically, for any state-action pair (s,a) and any vector  $QnR^{|S|}|_A|$ , we define

(4.3) 
$$T_{u,h}(Q)(s,a) = r(s,a) + a \quad E \quad \max_{s^e m P (t | s,a) \ pn \ a \ (A)} Q(s^e), p - uh_{s^e}(p) .$$

It is worth noting that this definition shares similarity with the regularized Bellman operator proposed in [21], where the operator defined there is targeted at  $V_u$ , while ours is defined w.r.t.  $Q_u$ .

The importance of this generalized Bellman operator is two-fold: it enjoys a desired contraction property, and its fixed point corresponds to the optimal regularized Q-function. These are generalizations of the properties for the classical Bellman operator, and are formally stated in the following lemma, whose proof can be found in [59].

Lemma 4.2 (properties of the generalized Bellman operator). For any u > 0, the operator  $T_{u,h}$  defined in (4.3) satisfies the following properties:

t  $T_{u,h}$  is a contraction operator w.r.t. the  $I_y$  norm; namely, for any  $Q_1,Q_2$  n R  $^{|S|}I_A^{|}$ , one has

$$| T_{u,h}(Q_1) - T_{u,h}(Q_2) |_{y} |_{q} |_{Q_1} - Q_2|_{y} .$$

t The optimal regularized Q-function  $Q_u^r$  is a fixed point of  $T_{u,h}$ , that is,

$$\mathsf{T}_{\mathsf{u},\mathsf{h}}(\mathsf{Q}_{\mathsf{u}}^{\mathsf{r}}) = \mathsf{Q}_{\mathsf{u}}^{\mathsf{r}}.$$

4.2. Proof of Theorem 3.1. Inspired by [12], our proof consists of (i) characterizing the dynamics of  $I_y$  errors and establishing a connection to a useful linear system with two variables and (ii) analyzing the dynamics of this linear system directly. In what follows, we elaborate on each of these steps.

Step 1: Error contraction and its connection to a linear system. With the assistance of the above preparations, we are ready to elucidate how to characterize the convergence behavior of  $|Q_l^r - Q_l^{(k+1)}|_y$ . Recalling the update rule of i  $^{(k+1)}$  (cf. (2.18c)), we can deduce that

$$Q_{i}^{r} - ui^{(k+1)} = a^{(Q_{i}^{r} - ui^{(k)})} + (1 - a)^{(Q_{i}^{r} - Q_{i}^{(k)})}$$

with a =  $\frac{1}{1+3}$ , thus indicating that

(4.6) 
$$|Q_{i}^{r} - ui^{(k+1)}|_{y}^{l} q a |Q_{i}^{r} - ui^{(k)}|_{y}^{l} + (1 - a)^{l} Q_{u}^{r} - Q_{i}^{(k)}|_{y}^{l}$$

Interestingly, there exists an intimate connection between  $|Q_{\iota}^{r}-Q_{\iota}^{(k+1)}|_{y}$  and  $|Q_{\iota}^{r}-ui^{(k+1)}|_{y}$  that allows us to bound the former term by the latter. This is stated in the following lemma, with the proof postponed to Appendix A.2.

Lemma 4.3. Set a =  $\frac{1}{1+a}$ -u. The iterates of Algorithm 2.1 satisfy

$$\left| Q_{u}^{r} - Q_{t}^{(k+1)} \right|_{y}^{l} = q a^{l} \left| Q^{r} - ui^{(k+1)} \right|_{y}^{l} + aa^{k+1} \left| Q^{(0)} - ui^{(0)} \right|_{y}^{l} .$$

The above inequalities (4.6) and (4.7) can be succinctly described via a useful linear system with two variables  $|Q_u^r - Q_t^{(k)}|_y$  and  $|Q_u^r - ui^{(k)}|_y$ , that is,

(4.8) 
$$x_{k+1} q A x_k + a a^{k+1} y$$
,

where

This forms the basis for proving Theorem 3.1.

Step 2: Analyzing the dynamics of the linear system (4.8). Before proceeding, we note that a linear system similar to (4.8) has been analyzed in [12, section 4.2.2]. We intend to apply the following properties that have been derived therein:

(4.10a) 
$$x_{k+1} q A^{k+1} \begin{bmatrix} x_0 + a(a^{-1}A - 1)^{-1}y \end{bmatrix},$$

(4.10b) 
$$a(a^{-1}A - I)^{-1}y = \begin{cases} 0 & 0 \\ |Q_t^{(0)} - ui^{(0)}|_{y_t} \end{cases}$$

(4.10b) 
$$a(a^{-1}A - I)^{-1}y = \begin{bmatrix} 0 & 1 \\ |Q_{L}^{(0)} - ui^{(0)}|_{y} \\ (1 - a)a + a \end{bmatrix}$$
(4.10c) 
$$A^{k+1} = \begin{bmatrix} (1 - a)a + a \\ 1 \end{bmatrix} \begin{bmatrix} 1 - a \\ 1 \end{bmatrix}$$

Substituting (4.10c) and (4.10b) into (4.10a) and rearranging terms, we reach

which taken together with the definition of  $x_{k+1}$  gives

Step 3: Controlling  $|i_u^r(s) - i_u^{(k+1)}(s)|_1$  and  $|V_u^r - V_u^{(k+1)}|_y$ . It remains to convert this result to an upper bound on  $|i_u^r(s) - i_u^{(k+1)}(s)|_1$  and  $|V_u^r - V_u^{(k+1)}|_y$ . By virtue of Lemma 2.1, there exist two vectors  $g_u^r(s)$  n I  $h_s(i_u^r(s))$ ,  $g^{(k+1)}(s)$  n I  $h_s(i_u^r(s))$  and two scalars  $c_s^r$ ,  $c_s^{(k+1)}$  n R that satisfy

$$\begin{cases} u^{-1}Q_u^r(s) - c_s^r 1 &= g_u^r(s), \\ i^{(k+1)}(s,t) - c_s^{(k+1)} 1 &= g^{(k+1)}(s). \end{cases}$$

It holds for all s n S that

where (i) results from  $h_s(i_u^{(k+1)}(s)) - h_s(i_u^r(s)) \neq eg^{(k+1)}(s)$ ,  $i_u^{(k+1)}(s) - i_u^r(s)$ e. Plugging (4.11) into (4.13) completes the proof for (3.1b).

When  $h_s$  is 1-strongly convex w.r.t. the  $l_1$  norm, we can invoke the strong monotonicity property of a strongly convex function [6, Theorem 5.24] to obtain

$$|i_{u}(s) - i^{(k+1)}(s)|_{1}^{2} q \stackrel{e}{i_{u}(s)} - i^{(k+1)}(s), g_{u}^{r}(s) - g^{(k+1)}(s) \stackrel{e}{\leftarrow}$$

$$= i_{u}^{r}(s) - i^{(k+1)}(s), g_{u}^{r}(s) + c_{s}^{r}1 - g^{(k+1)}(s) - c_{s}^{(k+1)}1$$

$$q |i_{u}(s) - i^{(k+1)}(s)|_{1}^{r} g^{r}(s) + c_{s}1 - g^{(k+1)}(s) - c_{s}^{(k+1)}1 |_{y}$$

$$= u^{-1} |i_{u}(s) - i^{(k+1)}(s)|_{1}^{r} Q_{u}^{r}(s) - ui^{(k+1)}(s)|_{y}^{r},$$

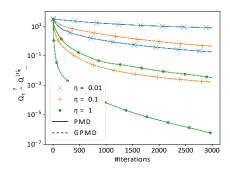
$$(4.14)$$

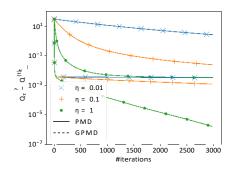
where the second line is valid since ei (s), 1e = ei (k+1)(s), 1e = 1. This taken together with (4.11) gives rise to the advertised bound

$$\left| i \left[ (s) - i \right]^{(k+1)}(s) \left| 1 - q \right|^{-1} \left| Q_{t}^{r}(s) - u i \right|^{(k+1)}(s) \left| 1 - q \right|^{-1} \left| Q_{t}^{r}(s) - u i \right|^{(k+1)}(s) \left| 1 - q \right|^{-1} \left| Q_{t}^{r}(s) - u i \right|^{-1}$$

- 5. Numerical experiments. In this section, we provide some simple numerical experiments to corroborate the effectiveness of the GPMD algorithm.
- 5.1. Tsallis entropy. While Shannon entropy is a popular choice of regularization, the discrepancy between the value function of the regularized MDP and the unregularized counterpart scales as  $O(\frac{u}{1-\log n}\log |A|)$ . In addition, the optimal policy un-der Shannon entropy regularization assigns positive mass to all actions and is hence nonsparse. To promote sparsity and obtain better control of the bias induced by regularization, the authors of [31, 32] proposed to employ the Tsallis entropy [49] as an alternative. To be precise, for any vector p n a (A), the associated Tsallis entropy is defined as

Tsallis<sub>q</sub>(p) = 
$$\frac{1}{q-1}$$
 1-  $\frac{m}{p(a)^{q}}$  =  $\frac{1}{q-1}$ E<sub>am p</sub> [ (p(a)) | q-1 ],





- (a) Tsallis entropy regularization
- (b) Log-barrier regularization

Fig. 1.  $|Q_t^r - Q_L^{(t)}|_y$  versus the iteration count for both PMD and GPMD, for multiple choices of the learning rate a. The left plot (a) is concerned with Tsallis entropy regularization, whereas the right plot (b) concerns log-barrier regularization used in our constrained RL example. The error curves are averaged over five independent runs.

where q > 0 is often referred to as the entropic index. When  $q \le 1$ , the Tsallis entropy reduces to the Shannon entropy.

We now evaluate numerically the performance of PMD and GPMD when applied to a randomly generated MDP with |S| = 200 and |A| = 50. Here, the transition probability kernel and the reward function are generated as follows. For each stateaction pair (s,a), we randomly select 20 states to form a set  $S_{s,a}$ , and set  $P(s^e|s,a) = 1/20$  if  $s^e$  n  $S_{s,a}$ , and 0 otherwise. The reward function is generated by r(s,a) m  $U_{s,a}$  t  $U_s$ , where  $U_{s,a}$  and  $U_s$  are independent uniform random variables over [0,1]. We shall set the regularizer as  $h_s(p) = -T$  sallis  $_2(p)$  for all s n s with a regularization parameter s u = 0.001. As can be seen from the numerical results displayed in Figure 1(a), GPMD enjoys a faster convergence rate compared to PMD.

5.2. Constrained RL. In reality, an agent with the sole aim of maximizing cumulative rewards might sometimes end up with unintended or even harmful behavior, due to, say, improper design of the reward function or nonperfect simulation of physical laws. Therefore, it is sometimes necessary to enforce proper constraints on the policy in order to prevent it from taking certain actions too frequently.

To simulate this problem, we first solve an MDP with |S| = 200 and |A| = 50, generated in the same way as in the previous subsection. We then pick 10 state-action pairs from the support of the optimal policy at random to form a set i . We can ensure that i  $|a| \le 1$  and  $|a| \le 1$  for all  $|a| \le 1$  and  $|a| \le 1$  and  $|a| \le 1$  for all  $|a| \le 1$  and  $|a| \le 1$  for all  $|a| \le 1$  and  $|a| \le 1$  for all  $|a| \le 1$  and  $|a| \le 1$  for all  $|a| \le 1$  and  $|a| \le 1$  for all  $|a| \ge 1$  fo

$$h_s(p) = \begin{cases} y & \text{if } (s,a) \text{ n i and } p(a) \neq i_{max}, \\ 0 & \text{if } (s,a) \text{ n i and } p(a) \neq i_{max}, \\ 0 & \text{otherwise.} \end{cases}$$

Numerical comparisons of PMD and GPMD when applied this problem are plotted in Figure 1(b). It is observed that PMD methods stall after reaching an error floor on the order of  $10^{-2}$ , while GPMD methods are able to converge to the optimal policy efficiently.

6. Discussion. The present paper has introduced a generalized framework of policy optimization tailored to regularized RL problems. We have proposed a

Generalized Policy Mirror Descent (GPMD) algorithm that achieves dimension-free linear convergence, which covers an entire range of learning rates and accommodates convex and possibly nonsmooth regularizers. Numerical experiments have been conducted to demonstrate the utility of the proposed GPMD algorithm. Our approach opens up a couple of future directions that are worthy of further exploration. For example, the current work restricts its attention to convex regularizers and tabular MDPs; it is of paramount interest to develop policy optimization algorithms when the regularizers are nonconvex and when sophisticated policy parameterization--including function approximation--is adopted. Understanding the sample complexities of the proposed algorithm--when the policies are evaluated using samples collected over an online trajectory--is crucial in sample-constrained scenarios and is left for future investigation. Furthermore, it might be worthwhile to extend the proposed algorithm to accommodate multi-agent RL, with a representative example being regularized multi-agent Markov games [14, 64, 11, 13].

Appendix A. Proof of key lemmas. In this section, we collect the proof of several key lemmas. Here and throughout, we use Ei [t] to denote the expectation over the randomness of the MDP induced by policy i. We shall follow the notation convention in (4.1) throughout. In addition, to further simplify notation, we shall abuse the notation by letting

(A.1a) 
$$D_{h_{s}}(\dot{e}, i; i) := D_{h_{s}}(\dot{e}(t|s), i(t|s); i(s,t)),$$
(A.1b) 
$$D_{h_{s}}(p, i; i) := D_{h_{s}}(p, i(t|s); i(s,t)),$$
(A.1c) 
$$D_{h_{s}}(i, p; i) := D_{h_{s}}(i(t|s), p; i(s,t))$$

(A.1b) 
$$D_{h_s}(p,i;i) := D_{h_s}(p,i;t|s);i(s,t),$$

(A.1c) 
$$D_{h_s}(i,p;i) := D_{h_s}(i,t|s),p;i(s,t)$$

for any policy i and  $\dot{\mathbf{e}}$  and any pn a (A), whenever it is clear from the context.

A.1. Proof of Lemma 2.1 We start by relaxing the probability simplex con- $_{an\ A}$  p(a) = 1 as straint (i.e., p n a (A )) in (2.18a) with a simpler linear constraint follows:

(A.2) minimize<sub>pn R | A |</sub> 
$$-a Q_t^{(k)}(s)$$
,  $p_e + auh_s(p) + D_{h_s}(p, i^{(k)}; i^{(k)})$  subject to  $p(a) = 1$ .

To justify the validity of dropping the non-negative constraint, we note that for any p obeying p(a) < 0 for some a n A, our assumption on  $h_s$  (see Assumption 1) leads to  $h_s(p) = y$ , which cannot possibly be the optimal solution. This confirms the equivalence between (2.18a) and (A.2).

Observe that the Lagrangian w.r.t. (A.2) is given by

$$\begin{array}{l} L_s \left(p, a_s^{(k)}\right) = -a^e Q_u^{(k)}(s), p^e + auh_s(p) + h_s(p) - h_s^{(i)}(s) - p - i^{(k)}(s), i^{(k)}(s)^e \\ \\ + a^{f(k)} & p(a) - 1 \end{array} ,$$

where  $a_s^{(k)}$  n R denotes the Lagrange multiplier associated with the constraint  $_{an\ A}$  p(a) = 1. Given that i  $^{(k+1)}$ (s) is the solution to (2.18a) and hence (A.2), the optimality condition requires that

$$0 \text{ n I }_{pL} \left( p, a_{s}^{(k)} \right) \Big|_{p=i} \left( k+1 \right)(s) = -aQ_{u}^{(k)}(s) + (1+au)I h_{s}^{(k+1)}(s) - i^{(k)}(s) + a_{s}^{(k)}1.$$

Rearranging terms and making use of the construction (2.15), we are left with

$$i^{(k+1)}(s) - \frac{\hat{\epsilon}_s^{(k)}}{1+au}1 = \frac{1}{1+au} [aQ_u^{(k)}(s) + i^{(k)}(s) - a_s^{(k)}1] n |h_s|^{(k+1)}(s)^{(k+1)}$$

thus concluding the proof of the first claim (2.16).

We now turn to the second claim (2.17). In view of the property (4.5), we have

i 
$$_{u}^{r}(s) = \underset{pn \ a \ (A)}{\text{rank}} \left( \frac{2.17}{r}, \frac{117}{r}, \frac{117}{r}, \frac{118}{r}, \frac{118}{r$$

This optimization problem is equivalent to

which can be verified by repeating a similar argument for (A.2). The Lagrangian associated with (A.3) is

$$(s, p, a_s^r) = -e_{Q_u^r(s), p}^r + uh_s(p) + a_s^r = p(a) - 1,$$

where  $a_s^r$  n R denotes the Lagrange multiplier. Therefore, the first-order optimality condition requires that

$$0 \text{ n I }_{pL_{s}}(p,a_{s}^{r})|_{p=i_{u}(s)}^{r} = -Q_{u}(s) + \text{ ul } h_{s}(i_{u}(s)) + a_{s}(s)$$

which immediately finishes the proof.

A.2. Proof of Lemma 4.3. Recall that  $Q_{\iota}^{(k+1)} = Q_{\iota}^{(k+1)}$ . In view of the relation (2.8), one obtains

$$Q_{l}^{(k+1)}(s,a) = r(s,a) + a \mathop{E}_{\substack{s^{e} \, m \, P \, (t \mid s,a) \\ s^{e} \, m \, (t \mid s,a) }} [ V_{l}^{(k+1)}(s^{e})$$

$$= r(s,a) + a \mathop{E}_{\substack{s^{e} \, m \, (t \mid s,a) \\ s^{e} \, m \, (t \mid s,a) }} [ \mathop{E}_{\substack{(k+1) \, (s^{e}) \\ \epsilon \, Q_{l}^{(k+1)}(s^{e}) }} [ Q_{l}^{(k+1)}(s^{e},a^{e}) - uh_{s^{e}}(i^{(k+1)}(s^{e}))^{]} ]$$

$$= r(s,a) + a \mathop{E}_{\substack{s^{e} \, m \, P \, (t \mid s,a) \\ s^{e} \, m \, P \, (t \mid s,a) }} [ Q_{l}^{(k+1)}(s^{e}), i^{(k+1)}(s^{e})^{-1} - uh_{s^{e}}(i^{(k+1)}(s^{e}))^{]} ]$$

This combined with the fixed-point condition (4.5) allows us to derive

$$\begin{aligned} Q_{u}^{\Gamma}(s,a) - & Q_{u}^{(k+1)}(s,a) \\ &= & T_{u,h}(Q_{u}^{\Gamma})(s,a) \\ &- & r(s,a) + a & E \\ &s^{e\,m\,P\,(t\,|s,a)} & \begin{bmatrix} e \\ Q_{u}^{(k+1)}(s^{e}), i^{(k+1)}(s^{e}) & e^{-uh_{s}e^{-i(k+1)}(s^{e})} \end{bmatrix} \end{bmatrix}^{1} \\ &= & T_{u,h}(Q_{u}^{\Gamma})(s,a) \\ &- & r(s,a) + a & E \\ &s^{e\,m\,P\,(t\,|s,a)} & \begin{bmatrix} e \\ ui^{(k+1)}(s^{e}), i^{(k+1)}(s^{e}) & -uh_{s}e^{-i(k+1)}(s^{e}) \end{bmatrix}^{1} \\ &- & a & E \\ &s^{e\,m\,P\,(t\,|s,a),a^{e}} & i^{(k+1)}(s^{e}) & -ui^{(k+1)}(s^{e},a^{e}) \end{bmatrix}^{1} \end{aligned}$$

In what follows, we control each term on the right-hand side of (A.4) separately.

Step 1: Bounding the 1st term on the right-hand side of (A.4). Lemma 2.1 tells us that

$$i^{(k+1)}(s) - c_s^{(k+1)} 1 n I h_s(i^{(k+1)}(s))$$

for some scalar  $c_s^{(k+1)}$  n R. This important property allows one to derive

(A.5) 
$$0 n - i^{(k+1)}(s) + c_s^{(k+1)} 1 + I h_s^{(i+1)}(s) = I L_{k+1,s}^{(i+1)}(s); c_s^{(k+1)},$$

where

$$L_{k+1,s}(p;a) := \frac{e}{\frac{e}{1}} \frac{e^{(k+1)}(s)}{e^{(k+1)}(s)} \frac{e}{\frac{e}{1}} \frac{e^{(k+1)}(s)}{e^{(k+1)}(s)} + a 1^{p} p.$$

Recognizing that the function  $f_{k+1,s}(t)$  is convex in p, we can view  $L_{k+1,s}(p;a)$  as the Lagrangian of the following constrained convex problem with Lagrangian multiplier an R:

(A.6) minimize 
$$f_{k+1,s}(p) = -\frac{e}{i^{(k+1)}}(s), p^{e} + h_{s}(p)$$
.

The condition (A.5) can then be interpreted as the optimality condition w.r.t. the program (A.6) and  $i^{(k+1)}(s)$ , meaning that

$$f_{k+1,s}$$
 ( $i^{(k+1)}(s)$ ) =  $\min_{p:1^p p=1} f_{k+1,s}(p)$ 

or, equivalently,

(A.7) 
$$e^{i(k+1)}(s), i^{(k+1)}(s)^e - h_s(i^{(k+1)}(s)) = \max_{p:1^p p=1} e^{i(k+1)}(s), p^e - h_s(p).$$

In addition, for any vector p that does not obey p q 0, Assumption 1 im-plies that  $h_s(p) = y$ , and hence p cannot possibly be the optimal solution to  $\max_{p \in i} (k+1)(s)$ , pe -  $h_s(p)$ . This together with (A.7) essentially implies that

(A.8) 
$$e_{i^{(k+1)}(s), i^{(k+1)}(s)}^{e} - h_{s^{(k+1)}(s)}^{e} = \max_{pn \ a \ (A)}^{e} e_{i^{(k+1)}(s), p}^{e} - h_{s^{(p)}}^{e}.$$

As a consequence, we arrive at

(A.9)

where the last step results from the contraction property (4.4) in Lemma 4.2.

Step 2: Bounding the 2nd term on the right-hand side of (A.4). Recall that  $a = \frac{1}{1+a}$ . Invoking the monotonicity property in Lemma 4.1 and the update rule (2.18c), we obtain

$$\begin{array}{l} Q^{(k+1)}(s,a) - ui^{(k+1)}(s,a) \\ = a Q_u^{(k+1)}(s,a) - ui^{(k)}(s,a)^{\frac{1}{2}} + (1-a)^{\frac{1}{2}} Q_u^{(k+1)}(s,a) - Q_u^{(k)}(s,a)^{\frac{1}{2}} \\ = q a Q_u^{(k)}(s,a) - ui^{(k)}(s,a) \ . \end{array}$$

Repeating this lower bound argument then yields

$$Q_{u}^{(k+1)}(s,a) - ui^{(k+1)}(s,a) q a^{k+1} Q_{u}^{(0)}(s,a) - ui^{(0)}(s,a)$$

$$q - a^{k+1} | Q_{u}^{(0)}(s,a) - ui^{(0)} |_{y},$$

thus revealing that

Step 3: Putting all this together. Substituting (A.9) and (A.10) into (A.4) gives

$$(A.11) \quad 0 \neq Q_{t}^{r}(s,a) - Q_{t}^{(k+1)}(s,a) \neq a | Q_{t}^{r} - ui^{(k+1)}|_{y} + a^{k+1} | Q_{t}^{(0)} - ui^{(0)}|_{y}$$

for all (s,a) n S s A, thus concluding the proof.

Appendix B. Analysis for approximate GPMD (Theorem 3.2). The proof consists of three steps: (i) evaluating the performance difference between i  $^{(k)}$  and i  $^{(k+1)}$ , (ii) establishing a linear system to characterize the error dynamic, and (iii) analyzing this linear system to derive global convergence guarantees. We shall describe the details of each step in what follows. As before, we adopt the notational convention (A.1) whenever it is clear from the context.

B.1. Step 1: Bounding performance difference between consecutive iterates. When only approximate policy evaluation is available, we are no longer guaranteed to have pointwise monotonicity as in the case of Lemma 4.1. Fortunately, we are still able to establish an approximate version of Lemma 4.1, as stated below.

Lemma B.1 (performance improvement for approximate GPMD). For all s n S and all k q 0, we have

$$V_u^{(k+1)}(s) \neq V_u^{(k)}(s) - \frac{1+a}{1-a}n_{opt} - \frac{2}{1-a}n_{eval}.$$

In addition, if  $h_s$  is 1-strongly convex w.r.t. the  $l_1$  norm for all s n S, then one further has

$$V_u^{(k+1)}(s) \neq V_u^{(k)}(s) - \frac{3+a}{1-a} n_{opt} - \frac{\epsilon}{(2+au)(1-a)} n_{eva}^2$$

In words, while monotonicity is not guaranteed, this lemma precludes the possibility of  $V_u^{(k+1)}(s)$  being much smaller than  $V_u^{(k)}(s)$ , as long as both  $n_{\text{eval}}$  and  $n_{\text{opt}}$  are reasonably small.

B.1.1. Proof of Lemma B.1. We divide the proof into two cases based on whether  $h_s$  is convex or strongly convex.

The case when  $h_s$  is convex. Let  $\dot{e}^{(k+1)}$  be the exact solution of the following problem:

problem: 
$$\{ e_{(k+1)}(s) = \arg\min_{p \in A(A)} - e_{(k)}(s), p_{e} + uh_{s}(p) + \frac{1}{\epsilon} D_{h_{s}}(p, i^{(k)}; i_{t}^{(k)}) \}$$

With this auxiliary policy iterate  $\dot{\mathbf{e}}^{(k+1)}$  in mind, we start by decomposing  $V_{\iota}^{(k+1)}(s)$ - $V_{\iota}^{(k)}(s)$  into the following three parts:

$$\begin{array}{l} & (B.2) \\ & (1-a)^{\binom{1}{k+1}}(s) - V_{\iota}^{(k)}(s)^{\binom{1}{k+1}}(s) - i_{\iota}^{(k)}(s^{e}) - i_{\iota}^{(k)}(s^{e})_{e} - uh_{s}^{e} i_{\iota}^{(k+1)}(s^{e})^{\binom{1}{k+1}}($$

where the first identity arises from the performance difference lemma for the regularized setting [27]. To continue, we seek to control each part of (B.2) separately.

t Regarding the first term of (B.2), straightforward computation indicates that

$$(B.3) \\ eQ_{()}^{k}(s_{[,i]}^{e},i_{(s_{[+1]}^{e})}^{(k+1)}(s_{[+1]}^{e}) + uh_{s_{[+1]}^{e}}(s_{[+1]}^{e},i_{(s_{[+1]}^{e})}^{(k+1)},i_{(s_{[+1]}^{e})}^{(k+1)}) \\ = \frac{1}{\epsilon}(1+au)D_{h_{s_{[+1]}^{e}}}(i_{[+1]}^{(k)},i_{[+1]}^{(k)}) + D_{h_{s_{[+1]}^{e}}}(i_{[+1]}^{(k)},i_{[+1]}^{(k)})$$

for all sen S.

t As for the second term of (B.2), the definition of the oracle  $G_{s,n}$  (see Assumption 3) guarantees that

$$(B.4) \\ - {}^{e}Q_{t}^{(k)}(s^{e}), i^{(k+1)}(s^{e})_{e} + uh_{se}^{(i^{(k+1)}(s^{e}))} + \frac{1}{\epsilon}D_{h_{se}}^{(i^{(k+1)}, i^{(k)}; t^{(k)})} \\ q - {}^{e}Q_{t}^{(k)}(s^{e}), \dot{e}^{(k+1)}(s^{e})_{e} + uh_{se}^{(\dot{e}^{(k+1)}(s^{e}))} + \frac{1}{\epsilon}D_{h_{se}}^{(\dot{e}^{(k+1)}, i^{(k)}; t^{(k)})} \\ + n_{opt}$$

for any sen S. Rearranging terms, we are left with

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

In addition, we note that the term

$$D_{h_{s^e}}$$
 ( $\dot{e}^{(k+1)}$ ,  $i^{(k)}$ ;  $i^{(k)}$ )

appears in both (B.3) and (B.5), which can be canceled out when summing these two equalities. Specifically, adding (B.3) and (B.5) gives

$$\begin{array}{c} {}^{e} \underline{Q}_{t}^{(k)}(s^{e}), \dot{e}^{(k+1)}(s^{e}) - i^{(k)}(s^{e})}_{e} - uh_{se} \stackrel{(\dot{e}^{(k+1)}(s^{e})}{e}) + uh_{se} \stackrel{(\dot{i}^{(k)}(s^{e})}{i^{(k)}(s^{e})} \\ + [\underline{Q}^{(k)}(s^{e}), i^{(k+1)}(s^{e}) - \dot{e}^{(k+1)}(s^{e}) - uh_{se} i^{(k+1)}(s^{e}) + uh_{se}] \dot{e}^{(k+1)}(s^{e})}_{1} \\ q \frac{1}{\epsilon} (1 + au)D_{h_{se}} \stackrel{(\dot{i}^{(k)}, \dot{e}^{(k+1)}; i^{(k+1)}}{i^{(k)}, \dot{e}^{(k+1)}; i^{(k+1)}} + D_{h_{se}} \stackrel{(\dot{i}^{(k+1)}, \dot{i}e^{(k)}; i^{(k)})}{i^{(k)}; i^{(k)}}_{1} \\ + \frac{1}{\epsilon} D_{h_{se}} \stackrel{(\dot{i}^{(k+1)}, \dot{i}^{(k)}; i^{(k)})}{i^{(k)}; i^{(k)}}_{1} - D_{h_{se}} \stackrel{(\dot{i}^{(k+1)}, \dot{e}^{(k)}; i^{(k)})}{i^{(k)}; i^{(k)}}_{1} - n_{opt}. \end{array}$$

Substituting this into (B.2) and invoking the elementary inequality  $|ea,be|q|a|_1|b|_y$  thus lead to

(B.6)

$$\begin{array}{l} V_{u}^{(k+1)}(s) - V_{u}^{(k)}(s) \\ q \ \frac{1}{1-a} \sum_{s^{e} \, md_{s}^{(k+1)}}^{[} \left[ \frac{1}{a} \left[ (1+au) \, D_{h_{s}e}^{\phantom{(k+1)}} (i^{(k)}, \dot{e}^{(k+1)}; it^{(k+1)}) + D_{h_{s}e}^{\phantom{(k+1)}} (i^{(k+1)}, \dot{e}^{(k)}; it^{(k)}) \right] \\ + \frac{1}{1-a} \sum_{s^{e} \, md_{s}^{(k+1)}}^{[} \left[ \frac{1}{a} \left( D_{h_{s}e}^{\phantom{(k+1)}} (i^{(k+1)}, i^{(k)}; it^{(k)}) - D_{h_{s}e}^{\phantom{(k)}} (i^{(k+1)}, \dot{e}^{(k)}; it^{(k)}) \right] \right] \\ - \frac{n_{pt}}{1-a} - \frac{1}{1-a} \sum_{s^{e} \, md_{s}^{(k+1)}}^{[} \left[ Q_{t}^{(k)}(s^{e}) - Q_{t}^{(k)}(s^{e}) \right]_{y}^{[} i^{(k+1)}(s^{e}) - i^{(k)}(s^{e}) \right]_{1}^{[} , \end{array}$$

where the last line makes use of Assumption 2 and the fact that  $|i|^{(k+1)}(s)|_1 = |i|^{(k)}(s)|_1 = 1$ .

Following the discussion in Lemma 2.1, we can see that  $i^{t(k)}(s) - c_s^{(k)} 1$  n I  $h_s(\dot{e}^{(k)}(s))$  with some constant  $c_s^{(k)}$  for all k. This together with the convexity of  $h_s$  (see (2.13)) guarantees that

(B.7) 
$$D_{h_s}(i^{(k)}, e^{(k+1)}; i^{t(k+1)}) \neq 0$$
 and  $D_{h_s}(i^{(k+1)}, e^{(k)}; t^{(k)}) \neq 0$ 

for any s n S, thus implying that the first term of (B.6) is non-negative. It remains to control the second term in (B.6). Towards this, a little algebra gives

(B.8) 
$$D_{h_s} = \frac{1}{1+au} \frac{(k+1)}{1+au} \frac{(k+1)$$

Here, the first and the third lines follow from the definition (2.13) and the second inequality comes from the construction (3.7), whereas the last step invokes the definition of the oracle (3.5). Substitution of (B.7) and (B.8) into (B.6) gives

$$\begin{array}{lll} \text{(B.9)} & & V_u^{(k+1)}(s) - V_u^{(k)}(s) \\ & & q - \frac{1+a}{1-a} n_{\text{opt}}, \\ & & - \frac{1}{1-e} \sum_{s^e \, \text{md}_s^{(k+1)}} \left[ \left| \, Q_u^{(k)}(s^e) - \, Q_u^{(k)}(s^e) \, \right|_{\, Y} \, \left| \, i^{\, (k+1)}(s^e) - \, i^{\, (k)}(s^e) \, \right|_{\, 1} \\ & & q - \frac{1+a}{1-a} n_{\text{opt}} - \frac{2}{1-a} n_{\text{eval}}. \end{array}$$

The case when  $h_s$  is strongly convex. When  $h_{s^e}$  is 1-strongly convex w.r.t. the  $l_1$  norm, the objective function of subproblem (B.1) is  $\frac{1+a}{a}$ -strongly convex w.r.t. the  $l_1$  norm. Taking this together with the  $n_{opt}$ -approximation guarantee in Assumption 3, we can demonstrate that

(B.11) 
$$\frac{1+aq}{2a} \stackrel{[e]}{\stackrel{(k+1)}{=}} (s^e) - i^{(k+1)}(s^e) \stackrel{|_2}{|_1} q n_{opt}$$
 for all k q 0 and  $s^e$  n S.

Additionally, the strong convexity assumption also implies that

$$\begin{split} & D_{h_{s^e}} \stackrel{\text{$($i^{(k+1)}(s^e)$ - $\dot{\mathbf{e}}^{(k)}(s^e)$; $i^{(k)}(s^e)$}{}^{0} q \frac{1}{2} |\dot{\mathbf{e}}^{(k)}(s^e)$ - $i^{(k+1)}(s^e)|_{1}^{2} \\ & = \frac{1}{2} \stackrel{\text{$($i^{(k+1)}(s^e)$ - $i^{(k+1)}(s^e)$}{}^{0} |_{1}^{2} + |\dot{\mathbf{e}}^{(k)}(s^e)$ - $i^{(k)}(s^e)|_{1}^{2} - \frac{1}{2} |\dot{\mathbf{e}}^{(k)}(s^e)$ - $i^{(k)}(s^e)$ - $i^{(k)}(s^e)|_{1}^{2} \\ & = \frac{1}{2} \stackrel{\text{$($i^{(k)}(s^e)$ - $i^{(k+1)}(s^e)$}{}^{0} |_{1}^{2} + |\dot{\mathbf{e}}^{(k)}(s^e)$ - $i^{(k)}(s^e)$ - $i^{(k)}(s^e)$$

where the third line results from Young's inequality, and the final step follows from (B.11). We can develop a similar lower bound on  $D_{h_se}(i^{(k)},e^{(k+1)};i^{(k+1)})$  as well. Taken together, these lower bounds give

$$\begin{split} &\frac{1}{\epsilon} \begin{bmatrix} (1+au)D_{h_{s^e}} & (i^{(k)}(s^e), \dot{e}^{(k+1)}(s^e); i^{(k+1)}(s^e) \\ + D_{h_{s^e}} & (i^{(k+1)}(s^e), \dot{e}^{(k)}(s^e); i^{(k)}(s^e) \end{bmatrix} \end{bmatrix}^{1} \\ & q \frac{2+au}{a} \begin{pmatrix} 1/4 |i^{(k)}(s^e) - i^{(k+1)}(s^e)|_{1}^{1} - \frac{ar_{oft}}{1+au} \\ q \frac{2+au}{4a} |i^{(k)}(s^e) - i^{(k+1)}(s^e)|_{1}^{1}^{2} - 2n_{opt}. \end{split}$$

In addition, it is easily seen that

$$\begin{array}{c} - \left| Q_{u}^{(k)}(s^{e}) - Q_{u}^{(k)}(s^{e}) \right|_{y} \left| i^{(k+1)}(s^{e}) - i^{(k)}(s^{e}) \right|_{1} \\ q - \frac{1}{2} \left( \frac{2a}{2 + au} \right| Q_{u}^{(k)}(s^{e}) - Q_{u}^{(k)}(s^{e}) \right|_{y}^{2} + \frac{2 + u}{2aa} \left| i^{(k+1)}(s^{e}) - i^{(k)}(s^{e}) \right|_{1}^{2} \\ q - \frac{\epsilon}{2 + au} \eta_{e_{v}al}^{2} - \frac{2 + au}{4a} i^{(k+1)}(s^{e}) - i^{(k)}(s^{e}) \right|_{1}^{2} \end{array}$$

Combining the above two inequalities with (B.10), we arrive at the advertised bound

$$V_u^{(k+1)}(s) - V_u^{(k)}(s) q - \frac{3+a}{1-a} n_{opt} - \frac{\epsilon}{(2+au)(1-a)} n_{eval}^2$$

B.2. Step 2: Connecting the algorithm dynamic with a linear system. Now we are ready to discuss how to control  $|Q_i^r - Q_i^{(k)}|_y$ . In short, we intend to establish the connection among several intertwined quantities and identify a simple linear system that captures the algorithm dynamic.

Bounding  $|Q_i^r - ui^{(k+1)}|_y$ . From the definition of  $i^{(k+1)}$ in (3.7), we have

$$\begin{split} \left|Q_{l}^{r} - u_{l_{1}(k+1)}\right|_{y} \\ &= \left|a\left(Q_{u}^{r} - u_{l_{1}(k)}\right) + (1-a)\left(Q_{u}^{r} - Q_{u}^{(k)}\right) + (1-a)\left(Q_{u}^{(k)} - Q_{u}^{(k)}\right)\right|_{y}^{l} \\ &= \left|a\left(Q_{u}^{r} - u_{l_{1}(k)}\right) + (1-a)\left(Q_{u}^{r} - Q_{u}^{(k)}\right) + (1-a)\left(Q_{u}^{(k)} - Q_{u}^{(k)}\right)\right|_{y}^{l} \\ &= \left|q\left(Q_{u}^{r} - u_{l_{1}(k)}\right)\right|_{y}^{l} + (1-a)\left|Q_{u}^{r} - Q_{u}^{(k)}\right|_{y}^{l} + (1-a)\left|Q_{u}^{(k)} - Q_{u}^{(k)}\right|_{y}^{l} \\ &= \left|q\left(Q_{u}^{r} - u_{l_{1}(k)}\right)\right|_{y}^{l} + (1-a)\left|Q_{u}^{r} - Q_{u}^{(k)}\right|_{y}^{l} + (1-a)n_{eval}, \end{split}$$

where the last inequality is a consequence of Assumption 2.

Bounding -  $\min_{s,a}(Q_t^{(k+1)}(s,a) - ut^{(k+1)}(s,a))$ . Applying the definition in (3.7) once again, we obtain

$$\begin{array}{c} \left(Q_{\zeta}^{(k+1)}(s,a) - u^{\dagger(k+1)}(s,a)\right) & (\\ = -a Q_{\zeta}^{(k)}(s,a) - u^{\dagger(k)}(s,a) + (1-a) Q_{\zeta}^{(k)}(s,a) - Q_{\zeta}^{(k)}(s,a) \\ & + \left(Q_{\zeta}^{(k)}(s,a) - Q_{\zeta}^{(k+1)}(s,a)\right) \\ q - a Q_{\zeta}^{(k)}(s,a) - u^{\dagger(k)}(s,a) + (1-a+c_1)n_{\text{eval}} + c_2n_{\text{opt}}, \end{array}$$

where

$$(B.14) \begin{array}{c} c_1 = \begin{cases} \begin{cases} \frac{2a}{1-a} & \text{if $h_s$ is convex but not strongly convex,} \\ \frac{a \, n_{\text{eval}} \, a}{(2+a \, u)(1-a)} & \text{if $h_s$ is $1$-strongly convex w.r.t. the $I_1$ norm,} \\ c_2 = \frac{\frac{(a+1)a}{1-a}}{\frac{(a+3)a}{1-a}} & \text{if $h_s$ is convex but not strongly convex,} \\ \end{cases} \\ \text{if $h_s$ is $1$-strongly convex w.r.t. the $I_1$ norm.} \end{array}$$

Here, the last step of (B.13) follows from Assumption 2 as well as the following relation:

$$Q_{\iota}^{(k)}(s,a) - Q_{\iota}^{(k+1)}(s,a) = a \sum_{s^{e} \text{ mP } (t \mid s,a)}^{E} V_{\iota}^{(k)}(s^{e}) - V_{\iota}^{(k+1)}(s^{e}) \neq c_{1} n_{eval} + c_{2} n_{opt},$$

where we have made use of Lemma B.1. Taking the maximum over (s,a) on both sides of (B.13) yields

(B.15) 
$$\begin{array}{c} ( & ) \\ - \min_{s,a} \ Q_{\iota}^{(k+1)}(s,a) - \ u^{\frac{1}{\epsilon}(k+1)}(s,a) \\ q - a \min_{s,a} \ (Q_{\iota}^{(k)}(s,a) - \ u^{\frac{1}{\epsilon}(k)}(s,a) \\ \end{array} \right) + (1 - a + c_1)n_{\text{eval}} + c_2n_{\text{opt}}.$$

Bounding  $|Q_u^r - Q_u^{(k+1)}|_y$ . To begin with, let us decompose  $Q_u^r(s,a) - Q_u^{(k+1)}(s,a)$  into several parts. Invoking the relation (4.5) in Lemma 4.2 as well as the property (2.8), we reach

In wat follows, we control the three terms in (B.16) separately.

t To begin with, we repeat an argument similar to that for (A.9) to show that

$$\begin{split} &T_{u,h}(Q_u^r)(s,a)\\ &-r(s,a)+a \sum_{s^e mP(t|s,a)} [e_{ui_t^{(k+1)}}(s^e),\dot{e}^{(k+1)}(s^e)_e - uh_{s^e}\dot{\dot{e}}^{(k+1)}(s)]]\\ &=T_{u,h}(Q_u^r)(s,a)-T_{u,h}(i^{t(k+1)})(s,a)\\ &q\ a\left|Q_t^r-ui^{t(k+1)}\right|_{\gamma}^l. \end{split}$$

t The second term of (B.16) can be bounded by applying (B.8) with k replaced by k + 1:

t As for the third term of (B.16), taking the maximum over all (s,a) n S s A gives

$$Q_u^{(k+1)}(s^e,a^e) - ut^{(k+1)}(s^e,a^e) \ q - \min_{s,a} \ Q_u^{(k+1)}(s,a) - ut^{(k+1)}(s,a) \ .$$

Taken together, the above bounds and the decomposition (B.16) lead to

A linear system of interest. Combining (B.12), (B.15), and (B.17), we reach the following linear system:

(B.18) 
$$z_{k+1} q B z_k + b$$
,

where

This linear system of three variables captures how the estimation error progresses as the iteration count increases.

B.3. Step 3: Linear system analysis. In this step, we analyze the behavior of the linear system (B.18) derived above. Observe that the eigenvalues and respective eigenvectors of the matrix B are given by

Armed with these, we can decompose z<sub>0</sub> in terms of the eigenvectors of B as follows:

$$\begin{array}{l} (B.22 \Big\{ \\ z_0 \ q \ | Q_t^r - Q_t^{(0)} |_y \\ | Q_t^{(0)} - uit^{(0)} |_y \\ = \frac{1}{a + (1 - a)a} \Big[ (1 \ a \ | Q_t^r - Q_u^{(0)} |_y + a \ | Q_u^r \ |_y + a \ | Q_u^{(0)} |_y + a \ | Q_u^{(0)} - uit^{(0)} |_y \\ + | Q_t^r - uit^{(0)} |_y v_2 + e_z v_3 \\ q \ \frac{1}{a + (1 - a)a} \Big[ V_{Q_t^r} Q_u^{(0)} |_y + 2 \ | Q_u^r \ uit^{(0)} |_y \Big] v_1 \\ + | Q_t^r - uit^{(0)} |_y v_2 + e_z v_3, \end{array}$$

where  $e_z$  n R is some constant that does not affect our final result. Also, the vector b defined in (B.19) satisfies

$$\begin{bmatrix} a \\ a(2-2a+c_1)n_{eval}+a(1-&+c_2)n_{opt} \\ ]bq & (1-a)n_{eval}+(1-a)n_{opt} \\ & (1-a+c_1)n_{eval}+c_2n_{opt} \\ & = \begin{bmatrix} (2-2a+c_1)n_{eval}+(1-a+c_2)n_{opt} \\ 1-a+c_2)n_{opt} \end{bmatrix} \begin{bmatrix} (1-a+c_1)n_{eval}+c_2n_{opt} \\ 1-a+c_2)n_{opt} \end{bmatrix} v_1 + \begin{bmatrix} (1-a+c_1)n_{eval}+c_2n_{opt} \\ 1-a+c_2)n_{opt} \end{bmatrix} v_2.$$

Using the decomposition in (B.22) and (B.23) and applying the system relation (B.18) recursively, we can derive

Copyright © by SIAM. Unauthorized reproduction of this article is prohibited.

$$\begin{array}{l} q \; B^{k+1} z_0 + \displaystyle \frac{m^k}{b} \; B^{k-1} b \\ q \; B^{k+1} \\ t \; \displaystyle \frac{1}{a + (1 - a)a} \left[ \left| Q_u^r - Q_t^{(0)} \right|_y + 2a \left| Q_t^r - u^{\frac{1}{2}(0)} \right|_y \right] v_1 + \left| Q_t^r - u^{\frac{1}{2}(0)} \right|_y v_2 + e_z v_3 \\ + \displaystyle \frac{m^k}{b} \; B^{k-1} \left[ (2 - 2a + c_1) n_{\,\text{eval}} + (1 - a + c_2) n_{\,\text{opt}} \right] v_1 + \left[ (1 - a + c_1) n_{\,\text{eval}} + c_2 n_{\,\text{opt}} \right] v_2 \\ = \left[ a_1^k \left( \left| Q_u^r - Q_t^{(0)} \right|_y + 2a \left| Q_t^r - u^{\frac{1}{2}(0)} \right|_y \right) \\ + \frac{1 - a_1^k + 1}{1 - a_1} \left[ (2 - 2a + c_1) n_{\,\text{eval}} + (1 - a + c_2) n_{\,\text{opt}} \right] v_1 \\ = \left[ \left( \left| Q_u^r - Q_t^{(0)} \right|_y + \frac{1 - a_1^k + 1}{1 - a_1} \right] \\ + \left( \left| Q_u^r - u^{\frac{1}{2}(0)} \right|_y + \frac{1 - a_1^k + 1}{1 - a_1^2} \right] (1 - a + c_1) n_{\,\text{eval}} + c_2 n_{\,\text{opt}} v_2. \end{array}$$

Recognizing that the first two entries of v2 are nonpositive, we can discard the term involving v2 and obtain

where

$$C := \frac{a}{1-1} \left[ (2-2a+c_1)n_{\text{eval}} + (1-a+c_2)n_{\text{opt}} \right].$$

Making use of the fact that 1 - 
$$a_1 = (1 - a)(1 - a)$$
, we can conclude that 
$$\begin{bmatrix} ( & & ) & & \\ & & 1 & \\ & & 1 - a & 2 + \frac{c_1}{1 - a} & n_{eval} + & 1 + \frac{c_2}{1 - a} & n_{opt} \\ \end{bmatrix}$$

The above bound essentially says that

and

Turning to  $V_u^r(s)$  -  $V_u^{(k+1)}(s)$ , by an argument similar to that of (4.13), we have

$$\begin{split} &V_{u}^{\Gamma}(s) - \ V_{u}^{(k+1)}(s) \\ &= \ Q_{u}^{L}(s) - \ Q_{c}^{(k+1)}(s), i^{(k+1)}(s) \\ &+ \ u(h_{s}(i^{(k+1)}(s)) - \ h_{s}(i^{\Gamma}_{u}(s))) - \ e^{Q_{u}^{\Gamma}(s), i^{(k+1)}(s) - i^{\Gamma}_{u}(s)} e^{e} \\ &= \ Q_{u}^{L}(s) - \ Q_{c}^{(k+1)}(s), i^{(k+1)}(s) + \ uD_{h_{s}}(i^{(k+1)}, i^{\Gamma}_{u}; g_{u}^{\Gamma}) \\ &q \ |Q_{u}^{\Gamma} - \ Q_{c}^{(k+1)}|^{1} + \ uD_{h_{s}}(i^{(k+1)}, h^{\sim (k+1)}; i^{L}_{c}(s), i^{L}_{c}(s) - g_{c}^{\Gamma}(s)^{\epsilon}, \\ &+ \ uD_{h_{s}}(h^{\sim (k+1)}, i^{\Gamma}_{c}; g_{c}^{\Gamma}) + \ u^{\epsilon}(i^{(k+1)}(s) - h^{\sim (k+1)}(s), i^{L}_{c}(s) - g_{c}^{\Gamma}(s)^{\epsilon}, \end{split}$$

where the third step results from the standard three-point lemma. To control the second term, we rearrange terms in (B.4) and reach at

$$\begin{split} n_{opt} \, q \, - \, & \stackrel{e}{\Phi}_{t}^{(k)}(s), i^{(k+1)}(s)^{e} \, + \, uh_{s}^{\, (i^{(k+1)}(s)^{\, )}} \, + \, \frac{1}{a} D_{\, h_{s}}^{\, (i^{(k+1)}, i^{(k)}; i^{\dagger}(k)^{\, )}} \\ & + \, \stackrel{e}{\Phi}_{t}^{(k)}(s), e^{(k+1)}(s)^{e} \, - \, uh_{s}^{\, (i^{(k+1)}(s)^{\, )}} \, - \, \frac{1}{\epsilon} D_{\, h_{s}}^{\, (i^{(k+1)}, i^{(k)}; i^{\dagger}(k)^{\, )}} \\ & = \, \stackrel{e}{\Phi}_{t}^{(k)}(s), e^{(k+1)}(s) \, - \, i^{(k+1)}(s)^{e} \, + \, \frac{1 + \, au}{\epsilon} h_{s}(i^{\, (k+1)}(s)) \, - \, h_{s}(h^{\sim (k+1)}(s))^{\, )} \\ & + \, \frac{1e}{\epsilon} t^{(k)}(s), e^{(k+1)}(s) \, - \, i^{(k+1)}(s)^{e} \\ & = \, \frac{1 + \, au}{\epsilon} D_{\, h_{s}}(i^{\, (k+1)}, h^{\sim (k+1)}; i^{\, (k+1)}) \, . \end{split}$$

For the remaining terms, recall that  $it^{(k+1)}$  -  $c_s^{(k+1)} 1$  n |  $h_s(\dot{e}^{(k+1)}(s))$  with some constant  $c_s^{(k+1)}$ . So we have

$$\begin{split} &uD_{\,h_{\,s}}\,({\textstyle h^{\sim}}^{(k+1)},i\,{\textstyle i\,\lceil}\,;g^r_\iota\,) + u\,{\textstyle i\,}^{\,e}\,{}^{(k+1)}(s) - {\textstyle h^{\sim}}^{(k+1)}(s), {\textstyle i\,\rceil}^{(k+1)}(s) - g^r_\iota\,(s) \\ &= uh_s({\textstyle h^{\sim}}^{(k+1)}(s)) - uh_s(i\,{}^{\,\epsilon}(s)) - {\textstyle h^{\sim}}^{(k+1)}(s) - {\textstyle i\,}^{\,\epsilon}\,{}^{\,\epsilon}(s), Q^{\,\epsilon}(s) \\ &+ u\,{\textstyle i\,}^{(k+1)}(s) - {\textstyle h^{\sim}}^{(k+1)}(s), {\textstyle i\,}^{(k+1)}(s) - g^r_\iota\,(s) \\ &q\,{\textstyle i\,}^r_\iota(s) - {\textstyle h^{\sim}}^{(k+1)}(s), Q^r_\iota\,(s) - ui^{N_{\,\ell}(k+1)}(s) \\ &- {\textstyle i\,}^{(k+1)}(s) - {\textstyle h^{\sim}}^{(k+1)}(s), Q^r_\iota\,(s) - ui^{N_{\,\ell}(k+1)}(s) \\ &= {\textstyle i\,}^r_\iota(s) - {\textstyle i\,}^{(k+1)}(s), Q^r_\iota\,(s) - ui^{N_{\,\ell}(k+1)}(s) \\ &q\,2{\textstyle I\,}^r_\iota(s) - uN^{N_{\,\ell}(k+1)}(s){\textstyle I\,}^r_\iota\,. \end{split}$$

Taken together, we conclude that

Finally, plugging in the choices of  $c_1$  and  $c_2$  (cf. (B.14)), we have C q  $C_2$  when  $\{h_s\}$  is convex and C q  $C_3$  when  $\{h_s\}$  is 1-strongly convex w.r.t. the  $I_1$  norm. In addition, for the latter case, we can follow an argument similar to that of (4.14) to demonstrate that

$$\begin{vmatrix} i_{u}^{r}(s) - \dot{e}_{u}^{(k)}(s) \end{vmatrix}_{1}$$

$$q u^{-1}(1-a)a + a)^{k} \begin{vmatrix} Q_{u}^{r} - Q_{u}^{(0)} \end{vmatrix}_{y} + 2a |Q_{u}^{r} - ui^{(0)}|_{y}^{l} + u^{-1}C_{3},$$

which taken together with (B.11) gives

$$\begin{vmatrix} i r (s) - i {k \choose u} (s) \end{vmatrix}^1 q \begin{vmatrix} i r (s) - e {k \choose u} (s) \end{vmatrix}^1 + \begin{vmatrix} i {k \choose u} (s) - e {k \choose u} (s) \end{vmatrix}^1$$

$$q u^{-1} (1 - a) a + a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + \begin{vmatrix} i {k \choose u} (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} (s) \end{vmatrix}^1 + 2a \begin{vmatrix} 1 r (s) - e {k \choose u} ($$

This concludes the proof.

## REFERENCES

- [1] Y. Abbasi-Yadkori, P. Bartlett, K. Bhatia, N. Lazic, C. Szepesvari, and G. Weisz, Politex: Regret bounds for policy iteration using expert prediction, in Proceedings of the International Conference on Machine Learning, PMLR, 2019, pp. 3692-3702.
- [2] A. Agarwal, N. Jiang, S. M. Kakade, and W. Sun, Reinforcement Learning: Theory and Algorithms, tech. report, 2019.
- [3] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan, Optimality and approximation with policy gradient methods in Markov decision processes, in Proceedings of the Conference on Learning Theory, PMLR, 2020, pp. 64--66.
- [4] A. Agazzi and J. Lu, Global optimality of softmax policy gradient with single hidden layer neural networks in the mean-field regime, in Proceedings of the International Conference on Learning Representations, ICLR, 2021.
- [5] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Man\(\frac{1}{6}\), Concrete Problems in Al Safety, manuscript, 2016.
- [6] A. Beck, First-Order Methods in Optimization, SIAM, Philadelphia, 2017, https://doi.org/10.1137/1.9781611974997.
- [7] A. Beck and M. Teboulle, Mirror descent and nonlinear projected subgradient methods for convex optimization, Oper. Res. Lett., 31 (2003), pp. 167–175.
- [8] D. P. Bertsekas, Dynamic Programming and Optimal Control, 4th ed., Athena Scientific, Belmont, MA, 2017.
- [9] J. Bhandari and D. Russo, Global Optimality Guarantees for Policy Gradient Methods, preprint, https://arxiv.org/abs/1906.01786, 2019.
- [10] J. Bhandari and D. Russo, A Note on the Linear Convergence of Policy Gradient Methods, preprint, https://arxiv.org/abs/2007.11120v1, 2020.
- [11] S. Cen, F. Chen, and Y. Chi, Independent Natural Policy Gradient Methods for Potential Games: Finite-Time Global Convergence with Entropy Regularization, preprint, https://arxiv.org/abs/2204.05466, 2022.
- [12] S. Cen, C. Cheng, Y. Chen, Y. Wei, and Y. Chi, Fast global convergence of natural policy gradient methods with entropy regularization, Oper. Res., 70 (2022), pp. 2563-2578.
- [13] S. Cen, Y. Chi, S. S. Du, and L. Xiao, Faster Last-Iterate Convergence of Policy Optimization in Zero-Sum Markov Games, preprint, https://arxiv.org/abs/2210.01050, 2022.
- [14] S. Cen, Y. Wei, and Y. Chi, Fast policy extragradient methods for competitive games with entropy regularization, in Proceedings of the 35th Conference on Neural Information Processing Systems, 2022, pp. 27952--27964.
- [15] Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh, A Lyapunov-based approach to safe reinforcement learning, in Proceedings of the 32nd International Conference on Neural Information Processing Systems, 2019, pp. 8103-8112.
- [16] Y. Chow, O. Nachum, and M. Ghavamzadeh, Path consistency learning in Tsallis entropy regularized MDPs, in Proceedings of the International Conference on Machine Learning, PMLR, 2018, pp. 979-988.
- [17] B. Dai, A. Shaw, L. Li, L. Xiao, N. He, Z. Liu, J. Chen, and L. Song, SBEED: Conver-gent reinforcement learning with nonlinear function approximation, in Proceedings of the International Conference on Machine Learning, PMLR, 2018, pp. 1125-1134.
- [18] D. Ding, X. Wei, Z. Yang, Z. Wang, and M. Jovanovic, Provably eficient safe exploration via primal-dual policy optimization, in Proceedings of the International Conference on Artificial Intelligence and Statistics, PMLR, 2021, pp. 3304-3312.
- [19] Y. Efroni, S. Mannor, and M. Pirotta, Exploration-Exploitation in Constrained MDPs, preprint, https://arxiv.org/abs/2003.02189, 2020.
- [20] M. Fazel, R. Ge, S. Kakade, and M. Mesbahi, Global convergence of policy gradient meth-ods for the linear quadratic regulator, in Proceedings of the International Conference on Machine Learning, 2018, pp. 1467--1476.
- [21] M. Geist, B. Scherrer, and O. Pietquin, A theory of regularized Markov decision processes, in Proceedings of the International Conference on Machine Learning, 2019, pp. 2160-2169.
- [22] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor, in Proceedings of the International Conference on Machine Learning, PMLR, 2018, pp. 1861--1870.
- [23] B. Hao, N. Lazic, Y. Abbasi-Yadkori, P. Joulani, and C. Szepesvári, Adaptive approximate policy iteration, in Proceedings of the International Conference on Artificial Intelligence and Statistics, PMLR, 2021, pp. 523-531.
- [24] S. M. Kakade, A natural policy gradient, in Proceedings of the 14th International Conference on Neural Information Processing Systems, 2001, pp. 1531–1538.

- [25] S. Khodadadian, P. R. Jhunjhunwala, S. M. Varma, and S. T. Maguluri, On the Linear Convergence of Natural Policy Gradient Algorithm, preprint, https://arxiv.org/ abs/2105.01424, 2021.
- [26] K. C. Kiwiel, Proximal minimization methods with generalized Bregman functions, SIAM J. Control Optim., 35 (1997), pp. 1142--1168, https://doi.org/10.1137/S0363012995281742.
- [27] G. Lan, Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes, Math. Program., 198 (2023), pp. 1059--1106.
- [28] G. Lan, Z. Lu, and R. D. Monteiro, Primal-dual first-order methods with O(1/n) iteration-complexity for cone programming, Math. Program., 126 (2011), pp. 1--29.
- [29] G. Lan and Y. Zhou, An optimal randomized incremental gradient method, Math. Program., 171 (2018), pp. 167--215.
- [30] N. Lazic, D. Yin, Y. Abbasi-Yadkori, and C. Szepesvari, Improved regret bound and experience replay in regularized policy iteration, in Proceedings of the International Conference on Machine Learning, PMLR, 2021, pp. 6032-6042.
- [31] K. Lee, S. Choi, and S. Oh, Sparse Markov decision processes with causal sparse Tsallis entropy regularization for reinforcement learning, IEEE Robot. Autom. Lett., 3 (2018), pp. 1466-1473.
- [32] K. Lee, S. Kim, S. Lim, S. Choi, and S. Oh, Tsallis Reinforcement Learning: A Unified Framework for Maximum Entropy Reinforcement Learning, preprint, https://arxiv.org/ abs/1902.00137, 2019.
- [33] G. Li, Y. Wei, Y. Chi, and Y. Chen, Softmax policy gradient methods can take exponential time to converge, Math. Program., to appear.
- [34] B. Liu, Q. Cai, Z. Yang, and Z. Wang, Neural trust region/proximal policy optimization attains globally optimal policy, in Proceedings of the 33rd International Conference on Neural Information Processing Systems, 2020, pp. 10565-10576.
- [35] Y. Liu, K. Zhang, T. Basar, and W. Yin, An improved analysis of (variance-reduced) policy gradient and natural policy gradient methods, in Proceedings of the 34th International Conference on Neural Information Processing Systems, 2021.
- [36] J. Mei, Y. Gao, B. Dai, C. Szepesvari, and D. Schuurmans, Leveraging non-uniformity in first-order non-convex optimization, in Proceedings of the International Conference on Machine Learning, PMLR, 2021, pp. 7555-7564.
- [37] J. Mei, C. Xiao, B. Dai, L. Li, C. Szepes\\arin and D. Schuurmans, Escaping the gravitational pull of softmax, in Proceedings of the 34th International Conference on Neural Information Processing Systems, 2021.
- [38] J. Mei, C. Xiao, C. Szepesvari, and D. Schuurmans, On the global convergence rates of soft-max policy gradient methods, in Proceedings of the International Conference on Machine Learning, PMLR, 2020, pp. 6820-6829.
- [39] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, Human-level control through deep reinforcement learning, Nature, 518 (2015), pp. 529-533.
- [40] T. M. Moldovan and P. Abbeel, Safe Exploration in Markov Decision Processes, manuscript, 2012.
- [41] A. S. Nemirovsky and D. B. Yudin, Problem Complexity and Method Eficiency in Optimization, manuscript, 1983.
- [42] G. Neu, A. Jonsson, and V. G\u00e0mez, A Unified View of Entropy-Regularized Markov Decision Processes, preprint, https://arxiv.org/abs/1705.07798, 2017.
- [43] M. L. Puterman, Markov Decision Processes: Discrete Stochastic Dynamic Programming, John Wiley & Sons, New York, 2014.
- [44] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, Trust region policy optimization, in Proceedings of the International Conference on Machine Learning, PMLR, 2015, pp. 1889-1897.
- [45] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, Proximal Policy Optimization Algorithms, preprint, https://arxiv.org/abs/1707.06347, 2017.
- [46] L. Shani, Y. Efroni, and S. Mannor, Adaptive Trust Region Policy Optimization: Global Convergence and Faster Rates for Regularized MDPs, preprint, https://arxiv.org/abs/ 1909.02769, 2019.
- [47] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, Policy gradient methods for reinforcement learning with function approximation, in Proceedings of the 12th International Conference on Neural Information Processing Systems, 1999, pp. 1057–1063.

- [48] M. Tomar, L. Shani, Y. Efroni, and M. Ghavamzadeh, Mirror Descent Policy Optimization, preprint, https://arxiv.org/abs/2005.09814, 2020.
- [49] C. Tsallis, Possible generalization of Boltzmann-Gibbs statistics, J. Statist. Phys., 52 (1988), pp. 479--487.
- [50] N. Vieillard, T. Kozuno, B. Scherrer, O. Pietquin, R. Munos, and M. Geist, Leverage the average: An analysis of K L regularization in reinforcement learning, in Proceedings of the 34th International Conference on Neural Information Processing Systems, 2021.
- [51] L. Wang, Q. Cai, Z. Yang, and Z. Wang, Neural policy gradient methods: Global optimal-ity and rates of convergence, in Proceedings of the International Conference on Learning Representations, 2019.
- [52] W. Wang, J. Han, Z. Yang, and Z. Wang, Global convergence of policy gradient for linear-quadratic mean-field control/game in continuous time, in Proceedings of the International Conference on Machine Learning, PMLR, 2021, pp. 10772-10782.
- [53] R. J. Williams, Simple statistical gradient-following algorithms for connectionist reinforcement learning, Mach. Learn., 8 (1992), pp. 229-256.
- [54] R. J. Williams and J. Peng, Function optimization using connectionist reinforcement learning algorithms, Connect. Sci., 3 (1991), pp. 241--268.
- [55] L. Xiao, On the Convergence Rates of Policy Gradient Methods, preprint, https://arxiv. org/abs/2201.07443, 2022.
- [56] P. Xu, F. Gao, and Q. Gu, Sample efficient policy gradient methods with recursive variance reduction, in Proceedings of the International Conference on Learning Representations, 2019.
- [57] T. Xu, Y. Liang, and G. Lan, A Primal Approach to Constrained Policy Optimization: Global Optimality and Finite-Time Analysis, preprint, https://arxiv.org/abs/2011.05869v1, 2020.
- [58] M. Yu, Z. Yang, M. Kolar, and Z. Wang, Convergent policy optimization for safe reinforcement learning, in Proceedings of the 33rd International Conference on Neural Information Processing Systems, 2020, pp. 3127-3139.
- [59] W. Zhan, S. Cen, B. Huang, Y. Chen, J. D. Lee, and Y. Chi, Policy Mirror Descent for Regularized Reinforcement Learning: A Generalized Framework with Linear Convergence, preprint, https://arxiv.org/abs/2105.11066, 2021.
- [60] J. Zhang, J. Kim, B. O'Donoghue, and S. Boyd, Sample Eficient Reinforcement Learning with REINFORCE, preprint, https://arxiv.org/abs/2010.11364, 2020.
- [61] J. Zhang, A. Koppel, A. S. Bedi, C. Szepesva\(\frac{1}{3}\), and M. Wang, Variational policy gradi-ent method for reinforcement learning with general utilities, in Proceedings of the 34th International Conference on Neural Information Processing Systems, 2021, pp. 4572-4583.
- [62] J. Zhang, C. Ni, Z. Yu, C. Szepesvari, and M. Wang, On the convergence and sample eficiency of variance-reduced policy gradient method, in Proceedings of the 34th International Conference on Neural Information Processing Systems, 2021, pp. 2228–2240.
- [63] K. Zhang, B. Hu, and T. Bacar, Policy optimization for H<sub>2</sub> linear control with H<sub>y</sub> robustness guarantee: Implicit regularization and global convergence, SIAM J. Control Optim., 59 (2021), pp. 4081-4109, https://doi.org/10.1137/20M1347942.
- [64] Y. Zhao, Y. Tian, J. Lee, and S. Du, Provably efficient policy optimization for two-player zero-sum Markov games, in Proceedings of the International Conference on Artificial Intelligence and Statistics, PMLR, 2022, pp. 2736--2761.