

---

# Learning Value Equivalent Models in Policy Improvement Path for Efficient MBRL

---

**Shenao Zhang**  
Georgia Tech

**Zhaoran Wang**  
Northwestern University

**Tuo Zhao**  
Georgia Tech

## 1 Introduction

Model-Based Reinforcement Learning (MBRL) plays a central role in Reinforcement Learning (RL) and is well known for its sample efficiency. Conventional MBRL algorithms which acquire a transition-predictive forward model enable fast re-planning or sample generation without having to query the original environment. However, *objective mismatch* arises between training the forward dynamics model and the overall goal of RL that maximizes the average policy value.

To solve the misalignment between model learning and decision-making, Value-Aware Model Learning (VAML) framework [8, 7, 22] is proposed. With the intuition to disregard aspects of the environment which are not related to decision-making, VAML assigns different weights to the transitions according to the *optimal* state value, which, however, is not known beforehand. Despite the various approximations presented by Farahmand et al., *e.g.*, learning in an iterative manner [7] or with worst-case loss [8], we turn our attention to another closely related decision-aware principle called *Value Equivalence* (VE) [20, 19, 10, 9], which constructs an abstract MDP model such that learning in it is equivalent to learning in the real environment. It can be done by learning VE models with respect to a set of functions and policies so that they yield the same Bellman updates.

The increasing number of policies and functions poses more stringent constraints on the VE model, which leads to the shrink of the VE model class. Intuitively, the model that is value equivalent with *all* possible policies and functions accurately describe the real MDP. Proper VE [9] loosed this requirement by showing that models that are VE w.r.t. all deterministic policies (and their value functions) are sufficient for optimal planning. However, this still does not suggest an algorithmic framework or provide the efficiency guarantees due to the impracticality to enumerate all possible policies in an environment. Hence, it is important to investigate the smallest policy sets for characterizing the VE model and the efficiency of the resulting MBRL algorithms, which we study in this work.

We show that the long-term performance of VE-based planning is only affected by the model's ability to approximate policy values in its *policy-improvement path*, defined as a sequence of policies given by an iterative policy update algorithm. This result provides a principled guide for choosing the *minimum and necessary* set of policies to characterize a VE model and the asymptotic performance. Although the policy-improvement path is priori unknown, the underlying idea to give good value predictions along it motivates the perspective from online learning. Specifically, the policy set can be updated sequentially for future improvement path prediction based on past policies. However, this raises another problem regarding how good can we except the resulting algorithm to be.

We mainly study two widely adopted MBRL frameworks, namely sample-based policy search and model-based planning. For sample-based algorithms, the model-generated samples are used to perform model-free policy optimization. We show that when instantiated with natural policy gradient, the *global* convergence rate depends on the Rademacher complexity of the VE model loss, for which we design several online learners based on the policy-improvement path and show a  $T^{-\frac{1}{2}}$  rate of convergence. We also analyze the iterative model learning approaches in convex settings and provide several extensions inspired by FTPL and optimistic FTPL with predictors in nonconvex settings.

We also analyze another important usage of VE principle in model-based planning. With a case study of MuZero [19], which demonstrates the superiority of VE models in complex task, we provide the regret bound of its variant considered in [9]. The result suggests the benefit of an online model learner in the iterative training process of MuZero. Recent implementations and experiments in [9] also verify such designs by observing performance improvement when incorporating past policies as additional VE losses. Besides, we prove the approximate local optimality with an online model learner in more general VE-based planning algorithms.

## 2 Preliminaries

**Value-Aware Model Learning** is a framework that addresses the *objective mismatch* issue [13, 6]. The conventional approach for one-step ahead predictive model learning is to minimize the KL-divergence between the empirical data and the model, while the overall goal of RL is to maximizes the average policy value. To solve the misalignment between model learning and decision making, Farahmand et al. proposed VAML [8, 7, 22] that minimizes the discrepancy between the Bellman optimality operator induced by model  $\hat{f}$  and the real  $f^*$  on the *optimal* value function  $V$ :

$$l(\hat{f})(s, a) = |\langle f^*(\cdot|s, a) - \hat{f}(\cdot|s, a), V \rangle|. \quad (1)$$

However,  $V$  is unknown during training and different approaches are proposed to approximate objective (1). For example, the optimal value in [8] is replaced with the supremum over a function space to minimize the worst-case discrepancy. IterVAML [7] performs planning in an iterative manner, with the unknown value replaced by the estimation in each iteration. Recently, VaGram [22] provides several modifications towards a more practical algorithm.

**Value Equivalence** is another principle that supports decision-aware model learning, which states that two models are value equivalent (VE) w.r.t. a set of functions  $\mathbb{V}$  and a set of policies  $\Pi := \{\pi|\pi : \mathcal{S} \mapsto P(\mathcal{A})\}$  if they yield the same updates under corresponding Bellman operators [10, 9]. The class  $\mathcal{F}$  of models that are VE with  $f^*$  are expressed as

$$\mathcal{F}^1(\Pi, \mathbb{V}) = \{\hat{f} \in \mathcal{F} : \hat{T}_\pi v = \mathcal{T}_\pi v, \forall \pi \in \Pi, v \in \mathbb{V}\}, \quad (2)$$

where Bellman operator  $\mathcal{T}_\pi[v](s) = \mathbb{E}[r(s, a) + v(s')|f^*, \pi]$  and  $\hat{\mathcal{T}}_\pi[v](s) = \mathbb{E}[\hat{r}(s, a) + v(s')|\hat{f}, \pi]$ .

Define the set of *all* possible policies and functions as  $\Pi_{\text{all}} := \{\pi|\pi : \mathcal{S} \mapsto P(\mathcal{A})\}$  and  $\mathbb{V}_{\text{all}} := \{v|v : \mathcal{S} \mapsto \mathbb{R}\}$ . It is shown in [10] that  $\mathcal{F}^1(\Pi_{\text{all}}, \mathbb{V}_{\text{all}})$  either contains only the environment or is empty. It makes intuitive sense since as more value functions are augmented to be considered, the model will less compromise the performance, and eventually collapse to the real environment model.

Grimm et al. [9] further reduced the required functions for characterizing an effective model sufficient for planning. They show that the VE principle can be formulated w.r.t. (deterministic) policy value functions only. Specifically, define an order- $k$  VE model class as

$$\mathcal{F}^k(\Pi, \mathbb{V}) = \{\hat{f} \in \mathcal{F} : \hat{\mathcal{T}}_\pi^k v = \mathcal{T}_\pi^k v, \forall \pi \in \Pi, v \in \mathbb{V}\}, \quad (3)$$

where  $\mathcal{T}_\pi^k$  denote  $k$  applications of  $\mathcal{T}_\pi$ . When  $k \rightarrow \infty$ , the *Proper Value Equivalence*  $\hat{\mathcal{T}}_\pi^\infty v = \mathcal{T}_\pi^\infty v$  becomes  $\hat{V}^\pi = V^\pi$ . The need for selecting functions  $\mathbb{V}$  is removed and  $\mathcal{F}^k(\Pi, \mathbb{V})$  becomes  $\mathcal{F}^k(\Pi)$ . With  $\hat{\mathcal{T}}$  induced by  $\hat{f}$ , the order- $k$  VE loss is defined by

$$l_{\text{VE}}^{\hat{f}}(\pi) = \left\| V^\pi - \hat{\mathcal{T}}_\pi^k V^\pi \right\|. \quad (4)$$

We note one property given in [9] regarding the structure of the model class: for any  $k \in \mathbb{Z}^+$ ,  $\mathcal{F}^\infty(\Pi) = \bigcap_{\pi \in \Pi} \mathcal{F}^k(\{\pi\}, \{V^\pi\})$ . This provides an alternative way to describe the order- $\infty$  PVE model as the intersection of  $\mathcal{F}^k$  with respect to the singleton policies in  $\Pi$  and their respective value functions. We denote any VE model in  $\mathcal{F}^\infty(\Pi)$  as  $\hat{f}_{\Pi}$ .

**Sequential Rademacher Complexity** is an important concept for online learning problems. We adopt the notation from [16] to define the sequential Rademacher complexity  $\mathfrak{R}_T(\mathbb{L}; \mathbf{x}, \mathbf{y})$  of the loss function class  $\mathbb{L} = \{(x, y) \mapsto \mathcal{L}((x, y); \hat{f})\}$ . Denote a sequence of Rademacher random variables

$\epsilon = (\epsilon_1, \dots, \epsilon_T)$ . For any  $\bar{\mathcal{X}}$ -valued tree  $\mathbf{x}$  and any  $\mathcal{Y}$ -valued tree  $\mathbf{y}$ ,  $\mathfrak{R}_T(\mathbb{L}) = \sup_{\mathbf{x}, \mathbf{y}} \mathfrak{R}_T(\mathbb{L}; \mathbf{x}, \mathbf{y})$ , where

$$\mathfrak{R}_T(\mathbb{L}; \mathbf{x}, \mathbf{y}) \stackrel{\Delta}{=} \mathbb{E}_{\epsilon, \xi} \left[ \sup_{\mathcal{L} \in \mathbb{L}} \sum_{t=1}^T \epsilon_t \mathcal{L}((\mathbf{x}(\epsilon), \xi_t), \mathbf{y}(\epsilon)) \right]. \quad (5)$$

### 3 Value Equivalence in the Policy-Improvement Path

It is shown in [9] that the optimal policy in  $\hat{f}_{\mathbb{I}_{\text{all}}}$  is also an optimal policy in the environment. In practice, however, it is a stringent requirement to enumerate every  $\pi \in \mathbb{I}_{\text{all}}$  and learn the corresponding (proper) VE model. We now investigate if and how set  $\mathbb{I}_{\text{all}}$  can be effectively reduced to its subset  $\mathbb{I} \subset \mathbb{I}_{\text{all}}$  for VE to still succeed.

#### 3.1 Structure of Policy-Improvement Path

We show that the application of value-equivalence principle can be tailored as the learning of a model that allows for accurate value predictions of policies in an algorithm's *policy-improvement path*.

Specifically, we define the policy-improvement path  $\mathfrak{P}(\cdot) = \{\pi_1, \dots, \pi_t, \dots, \pi_T\}$  as a sequence of policies given by an algorithm, *e.g.*,  $\mathfrak{P}(\hat{f}_{\mathbb{I}})$  formed by policy updates  $\pi_{t+1} = \text{argmax}_{\pi} \hat{\mathcal{T}}_{\pi}^k V^{\pi_t}$ , where  $\hat{\mathcal{T}}$  is induced by  $\hat{f}_{\mathbb{I}}$ . We provide two properties regarding the structure of  $\mathfrak{P}(\hat{f}_{\mathbb{I}})$ .

Firstly, suppose w.l.o.g. that we are interested in obtaining a unique optimal policy  $\pi^*$  at the end of training. When  $\pi_T = \pi^*$ , all possible policy-improvement paths compose a tree-like structure in which the optimal policy is the root, the first level has all policies  $\pi_{T-1}$  that lead to  $\pi^*$  with one iteration of policy update, and so on. This indicates that an optimal VE model should always give accurate value predictions of  $V^{\pi^*}$ .

Secondly, it suffices to set  $\mathbb{I}$  such that  $\hat{f}_{\mathbb{I}}$  is value equivalent with all  $\pi \in \mathfrak{P}(\hat{f}_{\mathbb{I}})$ . The long-term performance of VE-based planning is only affected by the model's ability to approximate policy values in its policy-improvement path. We formalize these two properties with the following theorem.

**Theorem 1.** *For a value equivalent model  $\hat{f}_{\mathbb{I}}$  learned on  $\mathbb{I}$ , suppose the order- $k$  VE error in  $\mathfrak{P}(\hat{f}_{\mathbb{I}})$  is bounded by some  $\epsilon \geq 0$  for state distribution  $\mu$ , i.e.,*

$$\|V^{\pi} - \hat{\mathcal{T}}_{\pi}^k V^{\pi}\|_{\mu} \leq \epsilon, \forall \pi \in \mathfrak{P}(\hat{f}_{\mathbb{I}}).$$

*Then for any reference policy  $\bar{\pi}$ , the value gap in the real MDP between  $\bar{\pi}$  and the asymptotic policy  $\pi_{\infty}$  in  $\mathfrak{P}(\hat{f}_{\mathbb{I}})$  is bounded by*

$$\limsup_{T \rightarrow \infty} \mathbb{E}_{\mu} [V^{\bar{\pi}} - V^{\pi_T}] \leq \left( e_{\bar{\pi}} + \epsilon + 2\epsilon\gamma^k (I - \gamma^k \hat{\mathcal{P}}_{\pi_T}^k)^{-1} \hat{\mathcal{P}}_{\pi_T}^k \right) (I - \gamma^k \hat{\mathcal{P}}_{\bar{\pi}}^k)^{-1} \stackrel{\Delta}{=} z_1(e_{\bar{\pi}}) + z_2(\epsilon),$$

*where  $e_{\bar{\pi}} = \|V^{\bar{\pi}} - \hat{\mathcal{T}}_{\bar{\pi}}^k V^{\bar{\pi}}\|_{\mu}$  is the VE error of  $\hat{f}_{\mathbb{I}}$  w.r.t. the reference policy  $\bar{\pi}$  and  $\hat{\mathcal{P}}_{\pi}^k$  is the stochastic matrix of  $k$ -step transitions under policy  $\pi$  and model  $\hat{f}_{\mathbb{I}}$ .*

Theorem 1 shows the structure of policy set in VE learning and the positive prospects of replacing the large  $\mathbb{I}_{\text{all}}$  with its subset. Specifically, we have the lower bound for the asymptotic value  $\mathbb{E}[V^{\pi_{\infty}}] \geq \mathbb{E}[V^{\bar{\pi}}] - z_1(e_{\bar{\pi}}) - z_2(\epsilon)$ . The reference policy  $\bar{\pi}$  can be set freely, *e.g.*, to the optimal policy  $\pi^*$ , and obtain the sufficient condition for *global* asymptotic optimality: small  $\epsilon$  and small  $e_{\pi^*}$ . Due to the intractability of  $\pi^*$ , we now turn our attention to the asymptotic convergence. Since smaller  $\epsilon$  indicates smaller  $z_2(\epsilon)$ , the lower bound of  $\mathbb{E}[V^{\pi_{\infty}}]$  is improved by minimizing  $\epsilon$ . In other words, the long-term performance of VE-based planning can be evaluated by the model's ability to approximate policy values in its policy-improvement path. In particular if  $\epsilon = 0$ , then  $z_2(\epsilon) = 0$  and  $\mathbb{E}[V^{\pi_{\infty}}] \geq \mathbb{E}[V^{\pi^*}] - z_1(e_{\pi^*})$ .

Unfortunately, this does not directly suggest a training objective or support a practical algorithm, as the policy-improvement path is not known beforehand. We now analyze how the theorem inspires a tractable model learning objective.

### 3.2 Value Equivalent Model Learning

From Thm. 1, the lower bound of  $\mathbb{E}[V^{\pi_\infty}]$  is improved by minimizing the VE loss w.r.t.  $\mathfrak{P}(\hat{f}_{\Pi})$ . Therefore, at the iteration  $t$  of a specific policy-improvement path  $\mathfrak{P}(\hat{f}_{\Pi})$ , define a metric for  $\Pi$  as the average discounted VE error:  $l_{\Pi}(\pi_t) = \sum_{\pi_t \leq i \leq T \in \mathfrak{P}(\hat{f}_{\Pi})} \gamma^{i-t} l_{\text{VE}}^{\hat{f}_{\Pi}}(\pi_i)$ .

We note that the policy  $V^{\pi_t}$  in  $\mathfrak{P}(\hat{f}_{\Pi})$  is not a deterministic variable. Specifically, the value  $V^{\pi_t}$  that is used to minimize VE loss is estimated with samples from real trajectories, which brings randomness to policy update in stochastic systems. Besides, different initial policies also lead to distinct policy improvement paths. Denote under model  $\hat{f}_{\Pi}$  all possible policy-improvement paths as the union  $\bigcup \mathfrak{P}(\hat{f}_{\Pi})$  and the objective for  $\Pi$  is thus

$$\min_{\Pi} \mathbb{E}_{\pi \sim \bigcup \mathfrak{P}(\hat{f}_{\Pi})} [l_{\Pi}(\pi)]. \quad (6)$$

Clearly, the size of  $\bigcup \mathfrak{P}(\hat{f}_{\Pi})$  is in general much smaller than the entire policy set  $\Pi_{\text{all}}$ , which avoids the enumeration of all possible policies in an environment. However, such inclusion relation still cannot lead to a practical objective since it is impossible to specify  $\bigcup \mathfrak{P}(\hat{f}_{\Pi})$ . For this reason, we turn our attention to set  $\Pi$  that is learned concurrently with policy during training, instead of predefined before training. Consider with an augmentation strategy  $g$  (e.g., identical mapping), policy  $g(\pi_t)$  is added to  $\Pi$  at each iteration  $t$ . We provide the total order theorem regarding  $l_{\Pi}$ .

**Theorem 2.** *For  $\Pi$  and  $\Pi'$  with identical augmentation strategy  $g$  everywhere but iteration  $i$ , then either  $l_{\Pi} \preccurlyeq l_{\Pi'}$  or  $l_{\Pi} \succcurlyeq l_{\Pi'}$ , where  $\preccurlyeq$  and  $\succcurlyeq$  represent element-wise vector inequalities.*

With the total order property, minimizing the expectation in Equation (6) can be achieved by adding policy  $g(\pi_t)$  to  $\Pi$  at iteration  $t$  such that  $l_{\Pi}(\pi_t)$  is minimized, since if  $l_{\Pi}$  is not minimized at  $\pi_t$  then the minimum is also not achieved at other  $\pi \in \bigcup \mathfrak{P}(\hat{f}_{\Pi})$ .

$$g(\pi_t) = \operatorname{argmin}_{\pi} \sum_{\pi_t \leq i \leq T \in \mathfrak{P}(\hat{f}_{\Pi})} \gamma^{i-t} l_{\text{VE}}^{\hat{f}_{\Pi \cup \{\pi\}}}(\pi_i). \quad (7)$$

The basic idea to give good predictions along the priori unknown policy-improvement path motivates the analysis of VE model learning from the online learning perspective, where data becomes available in a sequential order and is used to update the best predictor for future data at each step. Below we will discuss how this intuition can guide the designs of policy sets  $\Pi$ . We will also present different VE algorithmic frameworks including sample-based policy search and model-based planning in Sec. 4 and 5, and provide positive theoretical guarantees for both frameworks.

## 4 Sample-Based Policy Search with VE Models

### 4.1 Convergence Rate of Natural Policy Gradient

In this section, we study a specific model-based policy search framework that generates samples from a learned model and performs policy optimization with the samples. In particular, we analyze natural policy gradient with value equivalent models and study the property discounted MDPs with function approximations.

Denote the state value and state-action value by unrolling model  $\hat{f}$  with policy  $\pi$  as  $\hat{V}^\pi$  and  $\hat{Q}^\pi$ , respectively. That is,  $\hat{V}^\pi(s_0) = \sum_{h=1}^{\infty} \gamma^h \hat{r}(\hat{s}_h, \hat{a}_h)$ , where  $\hat{s}_h \sim \hat{f}(\cdot | \hat{s}_{h-1}, \hat{a}_{h-1})$ ,  $\hat{a}_h \sim \pi(\cdot | \hat{s}_h)$ , and the objective is  $J(\pi) = \mathbb{E}_{s_0 \sim \rho} [V^\pi(s_0)]$ . For natural gradient  $\tilde{\nabla}_\theta J(\theta_t)$ , a Fisher information matrix is adopted for update  $\theta_{t+1} - \theta_t = \eta \tilde{\nabla}_\theta J(\theta_t)$ , where

$$\tilde{\nabla}_\theta J(\theta) = F(\theta)^{-1} \nabla_\theta J(\theta) = F(\theta)^{-1} \mathbb{E}_{s \sim d_{\mu_1}^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta} [A^{\pi_\theta}(s, a) \nabla_\theta \log \pi_\theta(a | s)]. \quad (8)$$

In practice,  $\hat{J}(\pi_{\theta_t})$  is estimated with  $B$  samples replacing the expectation.

**Theorem 3.** Assume  $V_{\pi_\theta}$  is  $L$ -smooth in  $\theta$  and for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$  that  $\log \pi_\theta(a|s)$  is a  $\iota$ -smooth function of  $\theta$ . For  $\eta \leq \frac{1}{L}$ ,

$$\min_{t \in [T]} J(\pi^*) - J(\pi_{\theta_t}) \leq \frac{\log |\mathcal{A}|}{T\eta(1-\gamma)} + \frac{\iota}{T(1-L\eta)(1-\gamma)} \left( \mathbb{E}[J(\pi_{\theta_T}) - J(\pi_{\theta_1})] + \eta \sum_{t=1}^T (b_t + \frac{v_t}{2}) \right),$$

where the gradient bias  $b_t$  and the upper bound  $v_t$  of gradient variance are defined by  $b_t = \|\nabla_\theta J(\pi_{\theta_t}) - \mathbb{E}[\nabla_\theta \hat{J}(\pi_{\theta_t})]\|_2$  and  $v_t = \mathbb{E} \left[ \left\| \nabla_\theta \hat{J}(\pi_{\theta_t}) - \mathbb{E}[\nabla_\theta \hat{J}(\pi_{\theta_t})] \right\|_2^2 \right]$ .

We analyze the gradient bias and variance terms below as their sums might still be linear in  $T$ .

## 4.2 Online VE Model Learning

For convenience, we assume a bounded gradient variance. We show in Appendix 6.7 that a wide range of RL problems can be contained with the assumptions held.

**Assumption 1.** There exists an absolute constant  $c \geq 0$  such that  $v_t \leq c^2/B$  for all  $t \in [T]$ .

We first focus on the likelihood ratio gradient estimator that leverages the score function to calculate gradient. The gradient bias can be represented by  $b_{t+1} = \|V^{\pi_t} - \hat{T}_{\pi_t}^\infty V^{\pi_t}\|$ , which corresponds to a proper VE loss [9] w.r.t. a single policy  $\pi_t$ . Different from pure model-free algorithms, in MBPG (Alg. 1) multiple PG updates are performed with model-generated data (for sample complexity reduction). In other words, the learned model  $\hat{f}_t$  continues to have an impact on the following  $n$  policy updates. Thus, simply setting  $\hat{f}_t = \operatorname{argmin}_{\hat{f}} b_t$  only guarantees a low-bias gradient update from  $\pi_{t-1}$  to  $\pi_t$ , leading to potentially large bias for the following  $n-1$  policy gradients.

Formally, in the convergence rate Theorem 3, we need the accumulate  $\sum_{t=1}^T b_t$  to be small. According to the PG update in Eq. (8), we define at any given state the regret for the model learner as

$$\text{REG}_T \triangleq \sum_{t=1}^T \mathcal{L}_t \left( \left( (s, a), \hat{V}^\pi \right); \hat{f}_{\lfloor \frac{t}{n} \rfloor n} \right) \triangleq \sum_{t=1}^T \left\| V^{\pi_t} - \hat{T}_{\pi_t}^\infty V^{\pi_t} \right\|, \quad (9)$$

where  $\hat{T}_{\pi_t}$  is induced by model  $\hat{f}_{\lfloor \frac{t}{n} \rfloor n}$ . The objective then becomes  $\min_{\hat{f}_1, \dots, \hat{f}_{\lfloor \frac{T}{n} \rfloor n}} \text{REG}_T$ .

For an online model learner  $\mathcal{R}$  giving predictions based on  $\mathcal{R}(\{x_i, y_i\}_{i=1}^{t-1})$ , where  $x$  and  $y$  are training inputs  $(s, a)$  and predictions  $\hat{V}^\pi$ , standard results in [16] showed the existence  $\mathcal{R}$  that guarantees  $\text{REG}_t \leq 2n\mathfrak{R}_T(\mathbb{L})$ . Typically, we have the Rademacher complexity  $\mathfrak{R}_T(\mathbb{L})$  bounded by  $\tilde{O}(\sqrt{T})$ , which is desired in the convergence rate theorem. Next, we discuss several instantiations of  $\mathcal{R}$ .

**General model learners.** Unlike  $\Pi$  that contains predefined policies, we denote  $\Pi_t$  as the policy set that changes during training. We first provide a learner in general settings which maintains  $\Pi_t = \{\pi_i\}_{i=1}^{t-1}$ . Define the model learning error at iteration  $t$  as  $\delta_t = \sum_{i=1}^{t-1} \mathbb{E}_{s \sim d^{\pi_i}} \left\| (\mathcal{T}_{\pi_t}^k - \hat{\mathcal{T}}_{\pi_t}^k) V^{\pi_i} \right\|$ , where  $\hat{\mathcal{T}}_{\pi_t}$  is induced by model  $\hat{f}_{\lfloor \frac{t}{n} \rfloor n}$ .

**Proposition 1.** With model learning error  $\delta = \max_{t \in [1, T]} \delta_t$ ,

$$\text{REG}_T \leq \left( \delta + \frac{1}{(1-\gamma)} \right) \sqrt{T \left( 1 + \frac{1}{(1-\gamma)^2} \right) \log \left( 1 + \frac{T}{(1-\gamma)^2} \right)}. \quad (10)$$

Suppose we have access to an offline optimization oracle that gives a finite constant  $\sigma$ , then regret  $\text{REG}_T \leq O(\sqrt{T \log T})$ .

Selecting past policies in the policy-improvement path as  $\Pi_t$  has shown its effectiveness in recent implementation [9], which further improves the state-of-the-art MuZero performance on Atari tasks.

**Online convex learning.** Despite the bound developed in general settings, for linear function approximations, the model learning process becomes an online convex optimization problem, where online gradient descent suffices to obtain  $O(\sqrt{T})$  regret. In this case, we can simply leverage *iterative*

model learning algorithms that set  $\Pi_t = \pi_{t-1}$  and minimize  $\delta'_t = \|V^{\pi_{t-1}} - \hat{\mathcal{T}}_{\pi_{t-1}}^\infty V^{\pi_{t-1}}\|$ . When the underlying transition of a high-dimensional MDP can be linearly modeled, we may also first convert the states to an abstract space for the convexity to be satisfied.

**Extensions in nonconvex settings.** Several extensions of Proposition 1 can be obtained in nonconvex settings, e.g., with nonlinear function approximation. In nonconvex online learning, a heuristic yet provable algorithm is *Follow the Perturbed Leader* (FTPL), which adds a perturbation noise to the loss:  $\delta'_t = \delta_t + \sigma_t$ . For non-oblivious loss which is not adversarially chosen, e.g., the VE model loss, we may expect it predictable from the past. Thus, minimizing  $\delta'_t = \delta_t + M_t + \sigma_t$  with a predictor  $M_t$  leads to optimistic algorithms (OFTPL) and tighter regret bound when  $M_t$  encodes useful information. For example, we may simply set  $M_t = \mathcal{L}_{t-1}$  to emphasise the policy similarities between consecutive iterations, or importance weighted past losses  $M_t = \frac{1}{t-1} \sum_{i=1}^{t-1} \alpha_i \mathcal{L}_i$  for fading memory statistics, or combined with additional VE losses given by random policies.

We show that FTPL and OFTPL achieve  $O(\sqrt{T})$  regret with access to an offline optimization oracle.

**Proposition 2.** Denote the  $(\alpha, \beta)$ -approximate optimization oracle of function  $g$  as  $\mathcal{O}_{\alpha, \beta}(g)$  in the sense that if  $x^* = \mathcal{O}_{\alpha, \beta}(g)$ , then  $g(x^*) - \langle \sigma, x^* \rangle \leq \inf_x f(x) - \langle \sigma, x \rangle + \alpha + \beta \|\sigma\|_1$ . Suppose the losses encountered by the model learner are  $l$ -Lipschitz w.r.t  $l_1$  norm and  $\sigma_t$  is an exponential distribution where each coordinate is sampled from  $\text{Exp}(\lambda_t)$ . For  $\lambda = O(T^{-\frac{1}{2}})$ ,  $\alpha = O(T^{-\frac{1}{2}})$ , and  $\beta = O(T^{-1})$ , we have for FTPL and OFTPL that  $\mathbb{E}[\text{REG}_T] \leq O(\sqrt{T})$ .

**Discussions.** All the above designs guarantee an  $\tilde{O}(T^{-\frac{1}{2}})$  rate of convergence in Theorem 3. However, value equivalence might not suffice to support broader problems of interest such as pathwise gradient estimator and PG without natural gradients. Instead, policy-aware or gradient-aware model learning [1, 5] should be analyzed since they directly learn a model for accurate gradient estimation, by imposing the score magnitude as an additional weight.

When implementation, calculating the order- $\infty$  PVE loss might be costly or even impractical. By the properties of order- $\infty$  and order- $k$  VE model class (see Preliminaries), we can replace PVE losses with order- $k$  VE losses by fitting an additional (non)parametric critic function and replace  $\hat{V}^\pi$  with the  $k$ -step bootstrapping. Besides, the memory cost for saving all past policies might be huge. Some approximations can be adopted, such as periodically storing the parameters as in [9], or performing policy mixtures with different learning rates applied to each NN head and use their average as the bootstrap target for all heads [3].

## 5 VE-Model-Based Planning

Another important usage of models is to perform model-based planning. Much attention has been paid to the planning with value equivalent models, such as Predictron, MuZero, VPN, and PVE.

### 5.1 Case Study of MuZero

MuZero [19] is a modern architecture that shows the power of value equivalence principle in model-based planning. The agent maintains a predictive model that outputs abstract state transitions, values, and rewards, then the action returned by MCTS (with upper confidence bound applied to trees) is executed in the real MDPs. The two stages are carried out iteratively.

Grimm et al. [9] demonstrate the close connection between the VE loss and the MuZero objective. Specifically, MuZero's per-state model loss can be expressed as:

$$l(s_h) = \sum_{i=0}^k (V_{h+i} - v(z_h^i))^2 + (r_{h+i} - \hat{r}(z_h^i))^2, \quad (11)$$

where  $V_{h+i} = r_{h+i} + \dots + r_{h+i+k'-1} + v(z_{h+i+k'}^0)$  is the  $k'$ -step bootstrapping.  $z_h^i$  is the abstract state after  $i$  step model rollout from the real state  $s_h$ .  $s_{h:H}$ ,  $a_{h:H}$ , and  $r_{h:H}$  are the real trajectory.

It is shown that the MuZero model loss upper bounds the order- $k$  VE loss, i.e.,  $\|V^{\pi_t} - \hat{\mathcal{T}}_{\pi_t}^k V^{\pi_t}\|_{d_\pi}^2 \leq c \mathbb{E}_{d_\pi}[l(s)]$ . Optimizing the MuZero model is equivalent to finding a VE model w.r.t.  $\Pi_t = \{\pi_{t-1}\}$ .

Since experiments in [9] show that even better results can be obtained by adding VE loss induced by past policies to MuZero, we investigate why  $\Pi_t = \{\pi_i\}_{i=1}^{t-1}$  is preferred over  $\Pi_t = \{\pi_{t-1}\}$ .

---

**Algorithm 1** Model-Based Policy Gradient

---

```

Initialize policy  $\pi_0$ , model  $\hat{f}$ ,  $t = 1$ 
for  $\lfloor \frac{T}{n} \rfloor n$  iterations do
    Execute policy  $\pi_{t-1}$  in the real MDP
    Update VE model  $\hat{f}_t$  (by minimizing  $\delta_t$  or  $\delta'_t$ )
    for  $n$  gradient updates do
        Calculate gradient  $\nabla_\theta \hat{J}(\pi_{t-1})$ 
        Perform PG and obtain updated  $\pi_t$ 
         $t \leftarrow t + 1$ 
    end for
end for

```

---

**Algorithm 2** Model-Based Planning

---

```

Initialize policy  $\pi_0$ , model  $\hat{f}$ ,  $t = 1$ 
for  $\lfloor \frac{T}{n} \rfloor n$  iterations do
    Execute policy  $\pi_{t-1}$  in the real MDP
    Update VE model  $\hat{f}_t$ 
    for  $n$  planning iterations do
        Perform model-based planning (e.g., MCTS)
         $t \leftarrow t + 1$ 
    end for
    Return  $\pi_t$ 
end for

```

---

UCB1 is an algorithm for multi-armed bandit that achieves logarithm regret with the number of actions taken. When applied to MCTS (UCT [12] or PUCT [17]), convergence to the optimal policy can be shown, if given access to the accurate payoffs. However, learning a unified model for value prediction is exactly the goal of MuZero, which results in iterative executions of tree search and model update. Thus, the best action searched by MCTS is closely related to the quality of the learned model. To derive the global convergence of MuZero, we also need the predicted model values to converge to values in the real MDP. Formally, for the behavior policy  $\pi_t$  returned by MCTS at iteration  $t$ , define the regret as  $\text{REG}_T = \sum_{t=1}^T \int_s \rho(s)(V^* - V^{\pi_t})(s)$ .

**Proposition 3.** Denote the number of MCTS runs with the learned model in every iteration as  $n$ . Suppose  $n = O(T)$ . Denote the Rademacher complexity for the MuZero loss (Eq. (11)) class  $L$  where  $l \in L$  as  $\mathfrak{R}_T$ . Then there exists a model learner  $\mathcal{R}$  for which the regret  $\text{REG}_T$  is bounded by

$$\text{REG}_T \leq O(\mathfrak{R}_T) + O(\sqrt{T \log T}).$$

We can choose online model learners for  $\mathfrak{R}_T$  to be bounded by  $\tilde{O}(\sqrt{T})$ . For example, define the model loss as  $\sum_{\pi \in \Pi_t} \mathbb{E}_{d_\pi}[l(s)]$  and  $\Pi_t = \{\pi_i\}_{i=1}^{t-1}$ .

## 5.2 Planning in General Settings

Note that the UCB algorithms used in MCTS and MuZero can only be applied to independent arms. For more general settings such as nonlinear models, recent works [4] have shown pessimistic results towards *global* convergence. Specifically, Dong et al. show that the complexity (e.g., Eluder dimension [18, 14]) of nonlinear models cannot be polynomially bounded even for one-layer neural networks. It is thus suggested to find an  $(\epsilon_g, \epsilon_h)$ -approximate local optimal policy  $\pi^*$  and use local regret instead as the evaluation metric, defined as

$$\text{REG}_{\epsilon_g, \epsilon_h}^{\text{LC}}(T) = \sum_{t=1}^T \left( \sup_{\pi^* \in \mathfrak{C}_{\epsilon_g, \epsilon_h}} V^{\pi^*} - V^{\pi_t} \right)$$

**Proposition 4.** Define the model loss  $l_t = \|V^{\pi_{t-1}} - \hat{V}^{\pi_{t-1}}\| + \|V^{\pi_t} - \hat{V}^{\pi_t}\|$  and  $\mathfrak{R}_T(L)$  the Rademacher complexity for class  $L$ , where  $l \in L$ . Under Lipschitz assumptions, we have for the  $(\epsilon, 6\sqrt{\zeta}\epsilon)$ -regret that

$$\text{REG}_{\epsilon, 6\sqrt{\zeta}\epsilon}^{\text{LC}}(T) \leq O(\sqrt{T \mathfrak{R}_T}).$$

Notably, the loss  $l_t$  contains the prediction error w.r.t. both  $\pi_{t-1}$  at the previous iteration and  $\pi_t$ , the *upcoming* policy. The intuition is to optimize the model so that monotonic policy improvement can be obtained to achieve the local optima within the convex hull. The simplest way to meet such condition is to give accurate predictions for successive policy updates. This motivates the usage of VE objectives with online model learners, such as FTL-style algorithms.

## 6 Proofs

### 6.1 Proof of Theorem 1

*Proof.* Define in the real MDP the value gap between policy  $\pi_t$  and  $\bar{\pi}$  as  $x_t = V^{\bar{\pi}} - V^{\pi_t}$ , the policy gain between successive policies as  $g_t = V^{\pi_{t+1}} - V^{\pi_t}$ .

We decompose and bound  $x_{t+1}$  as

$$\begin{aligned} x_{t+1} &= V^{\bar{\pi}} - V^{\pi_{t+1}} \\ &= V^{\bar{\pi}} - \hat{\mathcal{T}}_{\bar{\pi}}^k V^{\bar{\pi}} + \hat{\mathcal{T}}_{\bar{\pi}}^k V^{\bar{\pi}} - \hat{\mathcal{T}}_{\bar{\pi}}^k V^{\pi_t} + \hat{\mathcal{T}}_{\bar{\pi}_T}^k V^{\pi_t} - \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_t} \\ &\quad + \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_t} - \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_{t+1}} + \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_{t+1}} - V^{\pi_{t+1}}. \end{aligned} \quad (12)$$

Due to the VE error bound and the optimality of policy  $\pi_{t+1}$ , i.e.,  $\pi_{t+1} = \operatorname{argmax}_{\pi} \hat{\mathcal{T}}_{\pi}^k V^{\pi_t}$ , we have  $\hat{\mathcal{T}}_{\bar{\pi}_T}^k V^{\pi_t} - \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_t} \leq 0$  and

$$\mathbb{E}[x_{t+1}] \leq e_{\bar{\pi}} + \gamma^k \hat{\mathcal{P}}_{\bar{\pi}}^k \mathbb{E}[x_t] - \gamma^k \hat{\mathcal{P}}_{\pi_{t+1}}^k \mathbb{E}[g_t] + \epsilon, \quad (13)$$

where the expectation is taken over  $\mu$ .

For the policy gain  $g_t$ , we have

$$\begin{aligned} \mathbb{E}[g_t] &= \mathbb{E}[V^{\pi_{t+1}} - V^{\pi_t}] \\ &= \mathbb{E}[V^{\pi_{t+1}} - \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_{t+1}} + \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_{t+1}} - \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_t} \\ &\quad + \hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_t} - \hat{\mathcal{T}}_{\pi_t}^k V^{\pi_t} + \hat{\mathcal{T}}_{\pi_t}^k V^{\pi_t} - V^{\pi_t}] \\ &\geq -2\epsilon + \gamma^k \hat{\mathcal{P}}_{\pi_{t+1}}^k \mathbb{E}[g_t], \end{aligned} \quad (14)$$

where the inequality holds since  $\hat{\mathcal{T}}_{\pi_{t+1}}^k V^{\pi_t} \geq \hat{\mathcal{T}}_{\pi_t}^k V^{\pi_t}$ .

Thus,

$$-\mathbb{E}[g_t] \leq 2\epsilon(I - \gamma^k \hat{\mathcal{P}}_{\pi_{t+1}}^k)^{-1}. \quad (15)$$

Plugging into Equation (13), we obtain

$$\mathbb{E}[x_{t+1}] \leq e_{\bar{\pi}} + \epsilon + 2\epsilon\gamma^k(I - \gamma^k \hat{\mathcal{P}}_{\pi_{t+1}}^k)^{-1} \hat{\mathcal{P}}_{\pi_{t+1}}^k I + \gamma^k \hat{\mathcal{P}}_{\bar{\pi}}^k \mathbb{E}[x_t]. \quad (16)$$

By taking the limit superior component-wise, we have

$$\limsup_{T \rightarrow \infty} \mathbb{E}[x_T] \leq \epsilon \left( I + \frac{e_{\bar{\pi}}}{\epsilon} + 2\gamma^k(I - \gamma^k \hat{\mathcal{P}}_{\pi_T}^k)^{-1} \hat{\mathcal{P}}_{\pi_T}^k \right) (I - \gamma^k \hat{\mathcal{P}}_{\bar{\pi}}^k)^{-1}.$$

□

### 6.2 Proof of Theorem 2

*Proof.* We first have

$$l_{\mathbb{P}}(\pi_t) = \gamma l_{\mathbb{P}}(\pi_{t+1}) + l_{\text{VE}}^{\hat{f}_{\mathbb{P}}}(\pi_t).$$

Define  $x(\pi)$  as the decision-index function for deterministic policy  $\pi$  that returns the vector  $c \in \mathbb{R}^{|\mathcal{A}|^{|S|}}$ , where the element in  $c$  that corresponds to  $(s, a)$  is 1 if  $\pi(a|s) = 1$  and 0 otherwise.

Further define  $Y_{\mathbb{P}}$  as a matrix with size  $|\mathcal{A}|^{|S|} \times |\mathcal{A}|^{|S|}$  such that  $Y_{\mathbb{P}} = x(\pi_t)^\top x(\pi_{t+1})$ , where  $\pi_{t+1} = \operatorname{argmax}_{\pi} \hat{\mathcal{T}}_{\pi}^k V^{\pi_t}$  with  $\hat{\mathcal{T}}$  induced by  $\hat{f}_{\mathbb{P}}$  (or replacing argmax with policy gradient update steps).

For  $\mathbb{P}$  and  $\mathbb{P}'$ , we have

$$\begin{aligned} l_{\mathbb{P}'} - l_{\mathbb{P}} &= l_{\text{VE}}^{\hat{f}_{\mathbb{P}'}} - l_{\text{VE}}^{\hat{f}_{\mathbb{P}}} + \gamma Y_{\mathbb{P}'} l_{\mathbb{P}'} - \gamma Y_{\mathbb{P}} l_{\mathbb{P}} \\ &= l_{\text{VE}}^{\hat{f}_{\mathbb{P}'}} - l_{\text{VE}}^{\hat{f}_{\mathbb{P}}} + \gamma(Y_{\mathbb{P}'} - Y_{\mathbb{P}})l_{\mathbb{P}'} + \gamma Y_{\mathbb{P}}(l_{\mathbb{P}'} - l_{\mathbb{P}}) \\ &= (I - \gamma Y_{\mathbb{P}})^{-1}(l_{\text{VE}}^{\hat{f}_{\mathbb{P}'}} - l_{\text{VE}}^{\hat{f}_{\mathbb{P}}} + \gamma(Y_{\mathbb{P}'} - Y_{\mathbb{P}})l_{\mathbb{P}}) \end{aligned} \quad (17)$$

Based on the setting of  $\mathbb{I}$  and  $\mathbb{I}'$ , we know  $l_{\text{VE}}^{\hat{l}_{\mathbb{I}'}} - l_{\text{VE}}^{\hat{l}_{\mathbb{I}}}$  is zero on its first  $k$  elements, and  $Y_{\mathbb{I}'} - Y_{\mathbb{I}}$  is zero on its first  $k$  rows.

Hence, the right-hand side of Eq. (17) is the product of a matrix with a vector whose first  $k$  elements are 0. We write  $C_i^{\mathbb{I}}$  as the column of matrix  $(I - \gamma Y_{\mathbb{I}})^{-1}$  that corresponds to  $\pi_i$  (column number  $|\mathcal{A}|^{|S|}$  is equal to the total number of deterministic policies, and  $\pi_i$  is one of them whose index can be found accordingly). Therefore,

$$l_{\mathbb{I}'} = l_{\mathbb{I}} + \alpha C_i^{\mathbb{I}}, \text{ with } \alpha \in \mathbb{R}.$$

Since  $(I - \gamma Y_{\mathbb{I}})^{-1} = \sum_{j=0}^{\infty} (\gamma Y_{\mathbb{I}})^j$ , whose entries are all non-negative, the entries of  $C_i^{\mathbb{I}}$  are also non-negative. Thus, either  $l_{\mathbb{I}} \preccurlyeq l_{\mathbb{I}'}$  or  $l_{\mathbb{I}} \succcurlyeq l_{\mathbb{I}'}$ , depending on the sign of  $\alpha$ .  $\square$

### 6.3 Proof of Theorem 3

*Proof.* Denote  $\beta_t = \frac{\theta_{t+1} - \theta_t}{\eta} = \tilde{\nabla}_{\theta} \hat{J}(\theta_t)$ . By Lipschitz assumption we have

$$\begin{aligned} J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t}) &\geq \nabla_{\theta} J(\pi_{\theta_t})^\top (\theta_{t+1} - \theta_t) - \frac{L}{2} \|\theta_{t+1} - \theta_t\|_2^2 \\ &= \eta \nabla_{\theta} J(\pi_{\theta_t})^\top \beta_t - \frac{L\eta^2}{2} \|\beta_t\|_2^2. \end{aligned} \quad (18)$$

Rewrite the exact gradient  $\nabla_{\theta} J(\pi_{\theta_t})$  as

$$\nabla_{\theta} J(\pi_{\theta_t}) = \left( \nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\nabla_{\theta} \hat{J}(\pi_{\theta_t})] \right) - \left( \nabla_{\theta} \hat{J}(\pi_{\theta_t}) - \mathbb{E}[\nabla_{\theta} \hat{J}(\pi_{\theta_t})] \right) + \nabla_{\theta} \hat{J}(\pi_{\theta_t}).$$

Then we bound  $\nabla_{\theta} J(\pi_{\theta_t})^\top \beta_t$  in Eq. (18) by bounding the resulting three terms.

$$\left| \left( \nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\nabla_{\theta} \hat{J}(\pi_{\theta_t})] \right)^\top \beta_t \right| \leq \|\beta_t\|_2 \|\nabla_{\theta} J(\pi_{\theta_t}) - \mathbb{E}[\nabla_{\theta} \hat{J}(\pi_{\theta_t})]\|_2 = \|\beta_t\|_2 b_t \quad (19)$$

$$\left( \nabla_{\theta} \hat{J}(\pi_{\theta_t}) - \mathbb{E}[\nabla_{\theta} \hat{J}(\pi_{\theta_t})] \right)^\top \beta_t \leq \frac{\|\nabla_{\theta} \hat{J}(\pi_{\theta_t}) - \mathbb{E}[\nabla_{\theta} \hat{J}(\pi_{\theta_t})]\|_2^2}{2} + \frac{\|\beta_t\|_2^2}{2} \quad (20)$$

$$\nabla_{\theta} \hat{J}(\pi_{\theta_t})^\top \beta_t \geq \|\beta_t\|_2^2, \quad (21)$$

where Eq. (21) holds due to the fact that  $\left( \theta_{t+1} - (\theta_t + \eta \nabla_{\theta} \hat{J}(\theta_t)) \right)^\top (\theta_{t+1} - \theta_t) \leq 0$ .

Thus, we can bound Eq. (18) by

$$J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t}) \geq \eta \left( -\|\beta_t\|_2 b_t - \frac{\|\nabla_{\theta} \hat{J}(\pi_{\theta_t}) - \mathbb{E}[\nabla_{\theta} \hat{J}(\pi_{\theta_t})]\|_2^2}{2} + \frac{\|\beta_t\|_2^2}{2} \right) - \frac{L\eta^2}{2} \|\beta_t\|_2^2. \quad (22)$$

By taking expectation in Eq. (22), we have

$$\begin{aligned} \left( \frac{\eta}{2} - \frac{L\eta^2}{2} \right) \mathbb{E}[\|\beta_t\|_2^2] &\leq \mathbb{E}[J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t})] + \eta \mathbb{E}[\|\beta_t\|_2] b_t + \frac{\eta}{2} v_t \\ &\leq \mathbb{E}[J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t})] + b_t + \frac{\eta}{2} v_t \end{aligned} \quad (23)$$

We have for  $\eta \leq \frac{1}{L}$  that

$$\mathbb{E}[\|\beta_t\|_2^2] \leq 2(\eta - L\eta^2)^{-1} \left( \mathbb{E}[J(\pi_{\theta_{t+1}}) - J(\pi_{\theta_t})] + \eta b_t + \frac{\eta}{2} v_t \right). \quad (24)$$

By smoothness, we have

$$\log \frac{\pi_{t+1}(a|s)}{\pi_t(a|s)} - \nabla_{\theta} \log \pi_t(a|s)(\theta_{t+1} - \theta_t) \geq -\frac{\iota}{2} \|\theta_{t+1} - \theta_t\|_2^2. \quad (25)$$

Then we obtain by the definition of KL divergence that

$$\begin{aligned}
\mathbb{E}_{s \sim d^{\pi^*}} [D_{\text{KL}}(\pi^* || \pi_t) - D_{\text{KL}}(\pi^* || \pi_{t+1})] &= \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*} \left[ \log \frac{\pi_{t+1}(a|s)}{\pi_t(a|s)} \right] \\
&\geq \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*} \left[ \nabla_{\theta} \log \pi_t(a|s)(\theta_{t+1} - \theta_t) - \frac{\iota}{2} \|\theta_{t+1} - \theta_t\|_2^2 \right] \\
&= \mathbb{E}_{s \sim d^{\pi^*}} \mathbb{E}_{a \sim \pi^*} \left[ \eta A^{\pi_t}(s, a) - \frac{\iota \eta^2}{2} \|\beta_t\|_2^2 \right] \\
&= (1 - \gamma) \eta (J(\pi^*) - J(\pi_t)) - \frac{\iota \eta^2}{2} \mathbb{E}[\|\beta_t\|_2^2].
\end{aligned} \tag{26}$$

where the last equality follows the performance difference lemma [11].

Thus,

$$\begin{aligned}
&\min_{t \in [T]} J(\pi^*) - J(\pi_{\theta_t}) \\
&\leq \frac{1}{T} \sum_{t=1}^T (J(\pi^*) - J(\pi_{\theta_t})) \\
&\leq \frac{1}{T\eta(1-\gamma)} \sum_{t=1}^T \mathbb{E}_{s \sim d^{\pi^*}} [D_{\text{KL}}(\pi^* || \pi_t) - D_{\text{KL}}(\pi^* || \pi_{t+1})] + \frac{\iota \eta^2}{2} \mathbb{E}[\|\beta_t\|_2^2] \\
&\leq \frac{\mathbb{E}_{s \sim d^{\pi^*}} [D_{\text{KL}}(\pi^* || \pi_1)]}{T\eta(1-\gamma)} + \frac{\iota}{T(1-L\eta)(1-\gamma)} \left( \mathbb{E}[J(\pi_{\theta_T}) - J(\pi_{\theta_1})] + \eta \sum_{t=1}^T (b_t + \frac{v_t}{2}) \right) \\
&\leq \frac{\log |\mathcal{A}|}{T\eta(1-\gamma)} + \frac{\iota}{T(1-L\eta)(1-\gamma)} \left( \mathbb{E}[J(\pi_{\theta_T}) - J(\pi_{\theta_1})] + \eta \sum_{t=1}^T (b_t + \frac{v_t}{2}) \right).
\end{aligned}$$

□

*Proof.* First, we have

$$\begin{aligned}
\|\hat{s}_{h+1} - \mathbb{E}[\hat{s}_{h+1}]\|_2 &= \left\| \hat{f}(\hat{s}_h, \hat{a}_h) + \xi_h - \mathbb{E}[\hat{f}(\hat{s}_h, \hat{a}_h)] \right\|_2 \\
&\leq L_{\hat{f}} \sqrt{\|\hat{s}_h - \mathbb{E}[\hat{s}_h]\|_2^2 + 1} + \sigma_{\xi} \\
&\leq L_{\hat{f}} \|\hat{s}_h - \mathbb{E}[\hat{s}_h]\|_2 + L_{\hat{f}} + \sigma_{\xi}.
\end{aligned} \tag{27}$$

Since  $\|s_0 - \mathbb{E}[s_0]\|_2 \leq \|\mathcal{S}\|_2$ , we have for  $h \geq 2$  that

$$\|\hat{s}_h - \mathbb{E}[\hat{s}_h]\|_2 \leq (L_{\hat{f}} + \sigma_{\xi}) \frac{L_{\hat{f}}^{h-2} - 1}{L_{\hat{f}} - 1} + L_{\hat{f}}^{h-1} \|\mathcal{S}\|_2$$

Besides,

$$\begin{aligned}
\left\| \nabla_{\theta} \hat{J}(\pi_{\theta_t}) - \mathbb{E}[\nabla_{\theta} \hat{J}(\pi_{\theta_t})] \right\|_2 &= \left\| \sum_{h=0}^k \nabla \hat{r}(\hat{s}_h, \hat{a}_h) - \mathbb{E}[\nabla \hat{r}(\hat{s}_h, \hat{a}_h)] \right\|_2 \\
&\leq \sum_{h=0}^k L_{\nabla \hat{r}} (\|\hat{s}_h - \mathbb{E}[\hat{s}_h]\|_2 + 1) \\
&\leq \left( \frac{L_{\nabla \hat{r}} (L_{\hat{f}} + \sigma_{\xi})}{L_{\hat{f}} - 1} + L_{\nabla \hat{r}} \|\mathcal{S}\|_2 \right) \frac{2L_{\hat{f}} (L_{\hat{f}}^k - 1)}{L_{\hat{f}} - 1} + k L_{\nabla \hat{r}}.
\end{aligned} \tag{28}$$

□

## 6.4 Proof of Proposition 1

*Proof.*

$$\begin{aligned}
\left\| V^{\pi_t} - \hat{\mathcal{T}}_{\pi_t}^{\infty} V^{\pi_t} \right\|^2 &= \left\| (\mathcal{T}_{\pi_t}^{\infty} - \hat{\mathcal{T}}_{\pi_t}^{\infty}) V^{\pi_t} \right\|^2 \\
&\leq \|V^{\pi_t}\|_{\mathcal{V}_{t-1}^{-1}}^2 \left( \sum_{i=1}^{t-1} \mathbb{E}_{s \sim d^{\pi_i}} \left\| (\mathcal{T}_{\pi_t}^{\infty} - \hat{\mathcal{T}}_{\pi_t}^{\infty}) V^{\pi_i} \right\|^2 + \frac{1}{(1-\gamma)^2} \right) \\
&\leq \left( \delta_t^2 + \frac{1}{(1-\gamma)^2} \right) \|V^{\pi_t}\|_{\mathcal{V}_{t-1}^{-1}}^2
\end{aligned} \tag{29}$$

where  $\mathcal{V}_t = 1 + \sum_{i=1}^t \mathbb{E}_{s \sim d^{\pi_i}} [V^{\pi_i}(s)^2]$ . By writing it in a recursive form,

$$\begin{aligned}
\mathcal{V}_t &= \mathcal{V}_{t-1} + \mathbb{E}_{s \sim d^{\pi_t}} [V^{\pi_t}(s)^2] \\
&= \mathcal{V}_{t-1} \left( 1 + \mathcal{V}_{t-1}^{-1} \mathbb{E}_{s \sim d^{\pi_t}} [V^{\pi_t}(s)^2] \right)
\end{aligned} \tag{30}$$

Then we have

$$\begin{aligned}
\log \mathcal{V}_t &= \log \mathcal{V}_{t-1} + \log \left( 1 + \mathcal{V}_{t-1}^{-1} \mathbb{E}_{s \sim d^{\pi_t}} [V^{\pi_t}(s)^2] \right) \\
&\geq \log \mathcal{V}_{t-1} + \frac{\mathcal{V}_{t-1}^{-1} \mathbb{E}_{s \sim d^{\pi_t}} [V^{\pi_t}(s)^2]}{1 + \mathcal{V}_{t-1}^{-1} \mathbb{E}_{s \sim d^{\pi_t}} [V^{\pi_t}(s)^2]} \\
&\geq \log \mathcal{V}_{t-1} + \frac{\|V^{\pi_t}\|_{\mathcal{V}_{t-1}^{-1}}^2}{1 + \frac{1}{(1-\gamma)^2}}.
\end{aligned} \tag{31}$$

Thus,

$$\begin{aligned}
\sum_{t=1}^T \|V^{\pi_t}\|_{\mathcal{V}_{t-1}^{-1}} &\leq \sqrt{T} \sqrt{\sum_{t=1}^T \|V^{\pi_t}\|_{\mathcal{V}_{t-1}^{-1}}^2} \\
&\leq \sqrt{T} \sqrt{\sum_{t=1}^T \left( 1 + \frac{1}{(1-\gamma)^2} \right) (\log \mathcal{V}_t - \log \mathcal{V}_{t-1})} \\
&= \sqrt{T} \sqrt{\left( 1 + \frac{1}{(1-\gamma)^2} \right) \log \mathcal{V}_T} \\
&\leq \sqrt{T \left( 1 + \frac{1}{(1-\gamma)^2} \right) \log \left( 1 + \frac{T}{(1-\gamma)^2} \right)}
\end{aligned} \tag{32}$$

Plugging Eq. (32) into Eq. (29) gives the result.  $\square$

## 6.5 Proof of Proposition 2

*Proof.* For regret  $\text{REG}_T$ , we have

$$\begin{aligned}
\text{REG}_T &= \sum_{t=1}^T \mathcal{L}_t(\hat{f}_{\lfloor \frac{T}{n} \rfloor n}) \\
&= \sum_{t=1}^T \left[ \mathcal{L}_t(\hat{f}_{\lfloor \frac{T}{n} \rfloor n}) - \mathcal{L}_t(\hat{f}_{t+1}) \right] + \sum_{t=1}^T \mathcal{L}_t(\hat{f}_{t+1}) \\
&\leq l \sum_{t=1}^T \|\hat{f}_{\lfloor \frac{T}{n} \rfloor n} - \hat{f}_{t+1}\|_1 + \sum_{t=1}^T \mathcal{L}_t(\hat{f}_{t+1}) \\
&= l \sum_{j=1}^{\lfloor \frac{T}{n} \rfloor n} \sum_{i=1}^n \|\hat{f}_{jn} - \hat{f}_{jn+i}\|_1 + \sum_{t=1}^T \mathcal{L}_t(\hat{f}_{t+1}) \\
&\leq l \sum_{j=1}^{\lfloor \frac{T}{n} \rfloor n} \sum_{i=1}^n \|\hat{f}_{jn} - \hat{f}_{jn+1}\|_1 + \dots + \|\hat{f}_{jn+i-1} - \hat{f}_{jn+i}\|_1 + \sum_{t=1}^T \mathcal{L}_t(\hat{f}_{t+1}) \\
&= l \sum_{j=1}^{\lfloor \frac{T}{n} \rfloor n} \sum_{i=1}^n (n+1-i) \|\hat{f}_{jn+i-1} - \hat{f}_{jn+i}\|_1 + \sum_{t=1}^T \mathcal{L}_t(\hat{f}_{t+1}).
\end{aligned}$$

Note that we only optimize the model at iteration  $\lfloor \frac{T}{n} \rfloor n$ , and the virtual  $\hat{f}_t$  at other iteration  $t$  is not what we have during training. It is introduced for the proof only.

By applying induction, we have for the expected regret

$$\begin{aligned}
\mathbb{E} \left[ \sum_{t=1}^T \mathcal{L}_t(\hat{f}_{\lfloor \frac{T}{n} \rfloor n}) \right] &\leq l \sum_{t=1}^T \frac{n(n+1)}{2} \mathbb{E}[\|\hat{f}_t - \hat{f}_{t+1}\|_1] + \sum_{t=1}^T \mathbb{E}[\mathcal{L}_t(\hat{f}_{t+1})] \\
&\leq l \sum_{t=1}^T \frac{n(n+1)}{2} \mathbb{E}[\|\hat{f}_t - \hat{f}_{t+1}\|_1] + \frac{d(\beta T + D)}{\lambda} + \alpha T.
\end{aligned} \tag{33}$$

The stability term  $\mathbb{E}[\|\hat{f}_t - \hat{f}_{t+1}\|_1]$  remains to be bounded. From Lemma 5 and Lemma 6 in [21], we have

$$\mathbb{E}[\|\hat{f}_t - \hat{f}_{t+1}\|_1] \leq 125\lambda D l d^2 + \frac{d\beta}{20\lambda l} + 2d\beta + \frac{\alpha}{20l}.$$

Plugging the above bound in Equation (33) gives the bound of expected regret.

$$\mathbb{E} \left[ \sum_{t=1}^T \mathcal{L}_t(\hat{f}_{\lfloor \frac{T}{n} \rfloor n}) \right] \leq O \left( T \lambda D l^2 d^2 n^2 + dl\beta + \frac{d(\beta T + D)}{\lambda} + \alpha T \right).$$

For  $\lambda = O(T^{-\frac{1}{2}})$ ,  $\alpha = O(T^{-\frac{1}{2}})$ , and  $\beta = O(T^{-1})$ , we have the  $O(\sqrt{T})$  expect regret bound.  $\square$

## 6.6 Proof of Proposition 3

*Proof.* Denote  $\hat{a}_t = \operatorname{argmax}_a \hat{Q}(s, a)$  and  $a_t$  as the policy action returned by MCTS with PUCT at iteration  $t$ . At any state  $s$ , we have

$$\begin{aligned}
TQ(s, a^*) - \sum_{t=1}^T Q(s, a_t) &= TQ(s, a^*) - \sum_{t=1}^T \hat{Q}(s, \hat{a}_t) + \sum_{t=1}^T \hat{Q}(s, \hat{a}_t) - \sum_{t=1}^T \hat{Q}(s, a_t) \\
&\quad + \sum_{t=1}^T \hat{Q}(s, a_t) - \sum_{t=1}^T Q(s, a_t) \\
&\leq TQ(s, a^*) - T\hat{Q}(s, a^*) + \sum_{t=1}^T \hat{Q}(s, \hat{a}_t) - \sum_{t=1}^T \hat{Q}(s, a_t) \\
&\quad + \sum_{t=1}^T \hat{Q}(s, a_t) - \sum_{t=1}^T Q(s, a_t) \\
&\leq V^{\pi^*} - \hat{V}^{\pi^*} + \hat{V}^{\hat{\pi}} - \hat{V}^{\pi} + \hat{V}^{\pi} - V^{\pi} \\
&\leq TQ(s, a^*) - T\hat{Q}(s, a^*) + \sum_{t=1}^T \hat{Q}(s, a_t) - \sum_{t=1}^T Q(s, a_t) \\
&\quad + O(T\sqrt{\frac{\log n}{n}}) \\
&\leq O(\mathfrak{R}_T) + O(\sqrt{T \log T}),
\end{aligned} \tag{34}$$

where the first inequality holds due to the optimality of  $\hat{a}_t$  on  $\hat{Q}$ , the second inequality follows the results in PUCT [17]. In the last inequality, the first term is due to the relationship between online learning and sequential Rademacher complexity, the second term is because  $n = O(T)$  and  $O(T\sqrt{\frac{\log n}{n}}) \leq O(T\sqrt{\frac{\log T}{T}}) \leq O(\sqrt{T \log T})$ .  $\square$

## 6.7 Examples

In Theorem 3 it is assumed that  $V_{\pi_\theta}$  is  $L$ -smooth in  $\theta$ . We first provide a possible condition in Example 1 for this assumption to hold when the reward  $r$  and transition function  $f$  are both Lipschitz continuous and smooth. Such assumptions are also consider by many previous works [2, 23, 15].

**Example 1.** (Bastani, 2021, Lemma D.2). Denote  $L_h$  as the Lipschitz constant for function  $h$  and  $\bar{L}_h = \max\{L_h, L_{\nabla h, 1}\}$ . Then  $\nabla V_{\pi_\theta}$  is  $L$ -Lipschitz, where  $L = 44H^5\bar{L}_r\bar{L}_f^{4H}$ . Particularly,  $\nabla_\theta V_{\pi_\theta}$  is Lipschitz continuous in  $\theta$  with Lipschitz constant  $24H^5\bar{L}_r\bar{L}_f^{4H}$ .

We also give an example setting that satisfies Assumption 1.

**Example 2.** With function approximations, for a stochastic model with form  $\hat{s}_{h+1} = \hat{f}(\hat{s}_h, \hat{a}_h) + \xi_h$ , where  $\xi_h \sim p(\xi)$  and  $p(\xi)$  is  $\sigma_\xi$ -subgaussian. Assume a Lipschitz continuous  $\hat{f}$  and  $\nabla \hat{r}$ . For initial state distribution  $\|s_0 - \mathbb{E}[s_0]\|_2 \leq \|\mathcal{S}\|_2$  and  $\|a\|_2 \leq 1$ , we have that  $v_r \leq c_t^2/B$  where

$$c_t = \left( \frac{L_{\nabla \hat{r}}(L_{\hat{f}} + \sigma_\xi)}{L_{\hat{f}} - 1} + L_{\nabla \hat{r}}\|\mathcal{S}\|_2 \right) \frac{2L_{\hat{f}}(L_{\hat{f}}^k - 1)}{L_{\hat{f}} - 1} + kL_{\nabla \hat{r}}. \tag{35}$$

Setting  $c$  as the largest  $c_t$  satisfies Assumption 1.

## References

- [1] Romina Abachi, Mohammad Ghavamzadeh, and Amir-massoud Farahmand. Policy-aware model learning for policy gradient methods. *arXiv preprint arXiv:2003.00030*, 2020.
- [2] Osbert Bastani. Sample complexity of estimating the policy gradient for nearly deterministic dynamical systems. In *International Conference on Artificial Intelligence and Statistics*, pages 3858–3869. PMLR, 2020.
- [3] Will Dabney, André Barreto, Mark Rowland, Robert Dadashi, John Quan, Marc G Bellemare, and David Silver. The value-improvement path: Towards better representations for reinforcement learning. *arXiv preprint arXiv:2006.02243*, 2020.
- [4] Kefan Dong, Jiaqi Yang, and Tengyu Ma. Provable model-based nonlinear bandit and reinforcement learning: Shelve optimism, embrace virtual curvature. *Advances in Neural Information Processing Systems*, 34, 2021.
- [5] Pierluca D’Oro, Alberto Maria Metelli, Andrea Tirinzoni, Matteo Papini, and Marcello Restelli. Gradient-aware model-based policy search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3801–3808, 2020.
- [6] Benjamin Eysenbach, Alexander Khazatsky, Sergey Levine, and Ruslan Salakhutdinov. Mismatched no more: Joint model-policy optimization for model-based rl. *arXiv preprint arXiv:2110.02758*, 2021.
- [7] Amir-massoud Farahmand. Iterative value-aware model learning. In *NeurIPS*, pages 9090–9101, 2018.
- [8] Amir-massoud Farahmand, Andre Barreto, and Daniel Nikovski. Value-aware loss function for model-based reinforcement learning. In *Artificial Intelligence and Statistics*, pages 1486–1494. PMLR, 2017.
- [9] Christopher Grimm, André Barreto, Gregory Farquhar, David Silver, and Satinder Singh. Proper value equivalence. *arXiv preprint arXiv:2106.10316*, 2021.
- [10] Christopher Grimm, André Barreto, Satinder Singh, and David Silver. The value equivalence principle for model-based reinforcement learning. *arXiv preprint arXiv:2011.03506*, 2020.
- [11] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *In Proc. 19th International Conference on Machine Learning*. Citeseer, 2002.
- [12] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer, 2006.
- [13] Nathan Lambert, Brandon Amos, Omry Yadan, and Roberto Calandra. Objective mismatch in model-based reinforcement learning. *arXiv preprint arXiv:2002.04523*, 2020.
- [14] Ian Osband and Benjamin Van Roy. Model-based reinforcement learning and the eluder dimension. *Advances in Neural Information Processing Systems*, 27, 2014.
- [15] Matteo Pirotta, Marcello Restelli, and Luca Bascetta. Policy gradient in lipschitz markov decision processes. *Machine Learning*, 100(2):255–283, 2015.
- [16] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Stochastic, constrained, and smoothed adversaries. *Advances in neural information processing systems*, 24, 2011.
- [17] Christopher D Rosin. Multi-armed bandits with episode context. *Annals of Mathematics and Artificial Intelligence*, 61(3):203–230, 2011.
- [18] Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. *Advances in Neural Information Processing Systems*, 26, 2013.
- [19] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- [20] David Silver, Hado Hasselt, Matteo Hessel, Tom Schaul, Arthur Guez, Tim Harley, Gabriel Dulac-Arnold, David Reichert, Neil Rabinowitz, Andre Barreto, et al. The predictron: End-to-end learning and planning. In *International Conference on Machine Learning*, pages 3191–3199. PMLR, 2017.
- [21] Arun Sai Suggala and Praneeth Netrapalli. Online non-convex learning: Following the perturbed leader is optimal. In *Algorithmic Learning Theory*, pages 845–861. PMLR, 2020.

- [22] Claas A Voelcker, Victor Liao, Animesh Garg, and Amir-massoud Farahmand. Value gradient weighted model-based reinforcement learning. In *International Conference on Learning Representations*, 2021.
- [23] Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2019.