

# Representations for Stable Off-Policy Reinforcement Learning

Dibya Ghosh<sup>1</sup> Marc G. Bellemare<sup>1</sup>

## Abstract

Reinforcement learning with function approximation can be unstable and even divergent, especially when combined with off-policy learning and Bellman updates. In deep reinforcement learning, these issues have been dealt with empirically by adapting and regularizing the representation, in particular with auxiliary tasks. This suggests that representation learning may provide a means to guarantee stability. In this paper, we formally show that there are indeed nontrivial state representations under which the canonical TD algorithm is stable, even when learning off-policy. We analyze representation learning schemes that are based on the transition matrix of a policy, such as proto-value functions, along three axes: approximation error, stability, and ease of estimation. In the most general case, we show that a Schur basis provides convergence guarantees, but is difficult to estimate from samples. For a fixed reward function, we find that an orthogonal basis of the corresponding Krylov subspace is an even better choice. We conclude by empirically demonstrating that these stable representations can be learned using stochastic gradient descent, opening the door to improved techniques for representation learning with deep networks.

## 1. Introduction

Value function learning algorithms are known to demonstrate divergent behavior under the combination of bootstrapping, function approximation, and off-policy data, what Sutton & Barto (2018) call the “deadly triad” (see also van Hasselt et al., 2018). In reinforcement learning theory, it is well-established that methods such as Q-learning and TD(0) enjoy no general convergence guarantees under linear function approximation and off-policy data (Baird, 1995;

Tsitsiklis & Roy, 1996). Despite this potential for failure, Q-learning and other temporal-difference algorithms remain the methods of choice for learning value functions in practice due to their simplicity and scalability.

In deep reinforcement learning, instability has been mitigated empirically through the use of auxiliary tasks, which shape and regularize the representation that is learned by the neural network. Methods using auxiliary tasks concurrently optimize the value function loss and an auxiliary representation learning objective such as visual reconstruction of observation (Jaderberg et al., 2016), latent transition and reward prediction (Gelada et al., 2019), adversarial value functions (Bellemare et al., 2019), or inverse kinematics (Pathak et al., 2017). In robotics, distributional reinforcement learning (Bellemare et al., 2017) in particular has proven a surprisingly effective auxiliary task (Bodnar et al., 2019; Vecerik et al., 2019; Cabi et al., 2019). While the stability of such methods remains an empirical phenomenon, it suggests that a carefully chosen representation learning algorithm may provide a means towards formally guaranteed stability of value function learning.

In this paper, we seek procedures for discovering representations that guarantee the stability of TD(0), a canonical algorithm for estimating the value function of a policy. We analyze the expected dynamics of TD(0), with the aim of characterizing representations under which TD(0) is provably stable. Learning dynamics of temporal-difference methods have been studied in depth in the context of a fixed state representation (Tsitsiklis & Roy, 1996; Borkar & Meyn, 2000; Yu & Bertsekas, 2009; Maei et al., 2009; Dalal et al., 2017). We go one step further by considering this representation as a component that can actively be shaped, and study stability guarantees that emerge from various representation learning schemes.

We show that the stability of a state representation is affected by: 1) the space of value functions it can express, and 2) how it parameterizes this space. We find a tight connection between stability and the geometry of the transition matrix, enabling us to provide stability conditions for algorithms that learn features from the transition matrix of a policy (Dayan, 1993; Mahadevan & Maggioni, 2007; Wu et al., 2018; Behzadian et al., 2019) and rewards (Petric, 2007; Parr et al., 2007). Our analysis reveals that a number of

<sup>1</sup>Google Research. Correspondence to: Dibya Ghosh <dibya.ghosh@berkeley.edu>.

popular representation learning algorithms, including proto-value functions, generally lead to representations that are not stable, despite their appealing approximation characteristics.

As special cases of a more general framework, we study two classes of stable representations. **The first class consists of representations that are approximately invariant under the transition dynamics (Parr et al., 2008), while the second consists of representations that remain stable under reparameterization.** From this study, we find that stable representations can be obtained from common matrix decompositions and furthermore, as solutions of simple iterative optimization procedures. Empirically, we find that different procedures trade off learnability, stability, and approximation error. In the large data regime, the Schur decomposition and a variant of the Krylov basis (Petrik, 2007) emerge as reliable techniques for obtaining a stable representation.

We conclude by demonstrating that these techniques can be operationalized using stochastic gradient descent on losses. We show that the Schur decomposition arises from the task of predicting the expectation of one’s own features at the next time step, whereas a variant of the Krylov basis arises as from the task of predicting future expected rewards. This is particularly significant, as both of these auxiliary tasks have in fact been heuristically proposed in prior work (François-Lavet et al., 2018; Gelada et al., 2019). Our result confirms the validity of these auxiliary tasks, not only for improving approximation error but, more importantly, for taming the famed instabilities of off-policy learning.

## 2. Background

We consider a Markov decision process (MDP)  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \rho, \gamma)$  on a finite state space  $\mathcal{S}$  and finite action space  $\mathcal{A}$ . The state transition distribution is given by  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ , the reward function  $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , the initial state distribution  $\rho \in \Delta(\mathcal{S})$ , and the discount factor  $\gamma \in [0, 1)$ . We write  $\mathcal{H} = \mathcal{S} \times \mathcal{A}$  with  $|\mathcal{H}| = n$ , and treat real-valued functions of state and action as vectors in  $\mathbb{R}^n$ .

A stochastic policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$  induces a Markov chain on  $\mathcal{H}$  with transition matrix  $P^\pi \in \mathbb{R}^{n \times n}$ . The value function  $Q^\pi \in \mathbb{R}^n$  for a policy  $\pi$  is the expected return conditioned on the starting state-action pair,

$$Q^\pi(s_i, a_i) = \mathbb{E}_\pi \left[ \sum_{t \geq 0} \gamma^t r(s_t, a_t) \mid s_0 = s_i, a_0 = a_i \right].$$

The value function also satisfies Bellman’s equation; in vector notation (Puterman, 1994),

$$Q^\pi = r + \gamma P^\pi Q^\pi$$

from which we recover the concise  $Q^\pi = (I - \gamma P^\pi)^{-1} r$ .

### 2.1. Approximate Policy Evaluation

Approximate policy evaluation is the problem of estimating  $Q^\pi$  from a family of value functions  $\{Q_\theta\}_{\theta \in \mathbb{R}^d}$  given a distribution of transitions  $(s, a, r, s') \sim \xi(s, a)P(s'|s, a)$  (c.f. Bertsekas, 2011). We refer to  $\xi \in \Delta(\mathcal{H})$  as the *data distribution*, and define  $\Xi \in \mathbb{R}^{n \times n}$  a diagonal matrix with the elements of  $\xi$  on the diagonal. If the data distribution is the stationary distribution of  $P^\pi$ , the data is *on-policy* and *off-policy* otherwise. We equip  $\mathbb{R}^n$  with the inner product and norm that is induced by the data distribution:  $\langle v_1, v_2 \rangle_\Xi = v_1^\top \Xi v_2$ . Most concepts from Euclidean inner products extend to this setting; see Appendix A for a review.

We consider a two-stage procedure for estimating value functions (Levine et al., 2017; Chung et al., 2019; Bertsekas, 2018). We first learn a *representation*, a  $d$ -dimensional mapping  $\phi : \mathcal{H} \rightarrow \mathbb{R}^d$ , through an explicit representation learning step. After a representation is learned, approximate policy evaluation is performed with the family of value functions linear in the representation  $\phi$ :  $Q_\theta(s, a) = \theta^\top \phi(s, a)$ , where  $\theta \in \mathbb{R}^d$  is a vector of weights.

The representation corresponds to a matrix  $\Phi \in \mathbb{R}^{n \times d}$  whose rows are the vectors  $\phi(s, a)$  for different state-action pairs  $(s, a)$ . For clarity of presentation, we assume that  $\Phi$  has full rank. A representation is *orthogonal* if  $\Phi^\top \Xi \Phi = I$ ; these correspond to features which are normalized and uncorrelated. We write  $\text{Span}(\Phi)$  to denote the subspace of value functions expressible using  $\Phi$ , and denote  $\Pi$  the orthogonal projection operator onto  $\text{Span}(\Phi)$ , with closed form  $\Pi = \Phi(\Phi^\top \Xi \Phi)^{-1} \Phi^\top \Xi$ .

### 2.2. Temporal Difference Methods

TD fixed-point methods are a popular class of methods for approximate policy evaluation that attempt to find value functions that satisfy  $Q = \Pi T^\pi Q$  (Bradtke & Barto, 1996; Gordon, 1995; Maei et al., 2009; Dann et al., 2014). If  $\Pi T^\pi$  has a fixed-point, the solution is unique (Lagoudakis & Parr, 2003) and can be expressed as

$$\theta_{TD}^* = (\Phi^\top \Xi (I - \gamma P^\pi) \Phi)^{-1} \Phi^\top \Xi r.$$

We study TD(0), the canonical update rule to discover this fixed point. With a step size  $\eta > 0$  and transitions sampled  $(s, a, r, s', a') \sim \xi(s, a)P(s'|s, a)\pi(a'|s')$ , TD(0) takes the update

$$\theta_{k+1} = \theta_k - \eta \nabla Q_{\theta_k}(s, a) (Q_{\theta_k}(s, a) - (r + \gamma Q_{\theta_k}(s', a'))).$$

In matrix form, this corresponds to an expected update over all state-action pairs:

$$\theta_{k+1} = \theta_k - \eta (\Phi^\top \Xi (I - \gamma P^\pi) \Phi) \theta_k - \Phi^\top \Xi r. \quad (1)$$

With appropriately chosen decay of the step size, the stochastic update will converge if the expected update converges

(Benveniste et al., 1990; Tsitsiklis & Roy, 1996). However, these updates are not the gradient of any well-defined objective function except in special circumstances (Barnard, 1993; Ollivier, 2018), and hence do not inherit convergence properties from the classical optimization literature. The main aim of this paper is to provide conditions on the representation matrix  $\Phi$  under which the update is convergent. We are especially interested in schemes that are convergent independent of the data distribution  $\xi$ .

We will characterize the stability of TD(0) and a representation through the spectrum of relevant matrices. For a matrix  $A \in \mathbb{R}^{k \times k}$ , the spectrum is the set of eigenvalues of  $A$ , written as  $\text{Spec}(A) = \{\lambda_1, \dots, \lambda_k\} \subset \mathbb{C}$ . The spectral radius  $\rho(A)$  denotes the maximum magnitude of eigenvalues. Stochastic transition matrices  $P^\pi$  satisfy  $\rho(P^\pi) = 1$ . We consider a potentially nonsymmetric matrix  $A \in \mathbb{R}^{k \times k}$  to be positive definite if all non-zero vectors  $x \in \mathbb{R}^k$  satisfy  $\langle x, Ax \rangle > 0$ .

### 2.3. Representation Learning

In reinforcement learning, a large class of methods have focused on constructing a representation  $\Phi$  from the transition and reward functions, beginning perhaps with proto-value functions (Mahadevan & Maggioni, 2007). Involving  $P^\pi$  and  $r$  in the representation learning process is natural, since the value function  $Q^\pi$  is itself constructed from these two objects. As we shall later see, the stability criteria for these are also simple and coherent. Additionally, there is a large body of literature on the ease (or difficulty) with which these methods can be estimated from samples, and by proxy are amenable to gradient-descent schemes. Here we review the most common of these representation learning methods along with a few obvious extensions. Table 1 shows how their construction arises from different matrix operations on  $P^\pi$  and, in the case of the Krylov basis, of  $r$ .

**Laplacian Representations:** Proto-value functions (Mahadevan & Maggioni, 2007) capture the high-level structure of an environment, using the bottom eigenvectors of the normalized Laplacian of an undirected graph formed from environment transitions. This formalism extends to reversible Markov chains with on-policy data, but does not generalize to directional transitions, stochastic dynamics, and off-policy data. In the general setting, the Laplacian representation (Wu et al., 2018) uses the top eigenvectors of the symmetrized transition matrix (EigSymm). We demonstrate in Section 4.3 that when data is off-policy, modifying the representation to omit eigenvectors whose eigenvalues exceed a threshold can provide strong stability guarantees.

**Singular Vector Representations:** Representations using singular vectors have been well-studied in representation learning for RL, because they are expressive and often yield strong performance guarantees. Fast Feature Selection (Be-

REPRESENTATION	DECOMPOSITION
PROTO-VALUE FUNCTIONS <sup>1</sup>	EIG( $P^\pi$ )
LAPLACIAN (EIGSYMM)	EIG( $(P^\pi + \Xi^{-1} P^\pi \Xi)$ )
SAFE EIGSYMM <sup>2</sup>	EIG( $(P^\pi + \Xi^{-1} P^\pi \Xi)$ )
SVD	SVD( $P^\pi$ )
SVD OF SUCCESSOR REP.	SVD( $(I - \gamma P^\pi)^{-1}$ )
SCHUR	SCHUR( $P^\pi$ )
KRYLOV BASIS	$\{r, P^\pi r, \dots, (P^\pi)^{d-1} r\}$
ORTHOG KRYLOV BASIS	ORTHOG( $\mathcal{K}_d(P^\pi, r)$ )

Table 1. Representation learning algorithms that learn features from the transition matrix and rewards. EIG is the spectral eigendecomposition, SVD the singular value decomposition, SCHUR the Schur decomposition, and ORTHOG an arbitrary orthogonal basis.  
<sup>1</sup> Only defined for reversible Markov chains with on-policy data.  
<sup>2</sup> Discards a partial set of features (see Section 4.3).

hzadian et al., 2019) uses the top left singular vectors of the transition matrix as features. Similarly, Stachenfeld et al. (2014) and Machado et al. (2018) use the top left singular vectors of the successor representation (Dayan, 1993), a time-based representation which predicts future state visitations:  $\Psi = (I - \gamma P^\pi)^{-1}$ . We discover in Section 3.4 that the SVD objective of minimizing the norm of approximation error fails to preserve the spectral properties of transition matrices needed for stability, and can induce divergent behavior in TD(0). In contrast, we show that decompositions constrained to preserve the spectrum of the transition matrix, such as the Schur decomposition, guarantee stability and performance.

**Reward-Informed Methods:** If the reward structure of the problem is known apriori, a representation can focus its capacity on modelling future rewards and how they diffuse through the environment. Towards this goal, Petrik (2007) suggested the Krylov basis generated by  $P^\pi$  and  $r$  as features. Bellman Error Basis Functions (BEBFs) (Parr et al., 2007) iteratively builds a representation by adding the Bellman error for the best solution found so far as a new feature. Parr et al. (2008) show that under certain initial conditions for BEBFs, both representations span the Krylov subspace  $\mathcal{K}_d(P^\pi, r)$  generated by rewards. Although no general guarantees exist for arbitrary rewards, we discover that when rewards are easily predictable, orthogonal representations that span this Krylov subspace have stability guarantees.

## 3. Stability Analysis of Arbitrary Representations

To begin, we study the stability of TD(0) given an arbitrary representation. For conciseness, we call TD(0) the algorithm whose expected update is described by equation 1; this is an algorithm which may or may not be off-policy (according to  $\Xi$  and  $P^\pi$ ), and learns a linear approximation of the value function  $Q^\pi$  using features  $\Phi$ . The following formalizes our

notion of stability.

**Definition 3.1.** *TD(0) is **stable** if there is a step-size  $\eta > 0$  such that when taking updates according to equation 1 from any  $\theta_0 \in \mathbb{R}^d$ , we have  $\lim_{k \rightarrow \infty} \theta_k = \theta_{TD}^*$ .*

### 3.1. Learning Dynamics

For a sufficiently small step-size  $\eta$ , the discrete update of equation 1 behaves like the continuous-time dynamical system

$$\frac{\partial}{\partial t}(\theta_t - \theta_{TD}^*) = -A_\Phi(\theta_t - \theta_{TD}^*), \quad (2)$$

whose behaviour is driven by the *iteration matrix*

$$A_\Phi = \Phi^\top \Xi (I - \gamma P^\pi) \Phi.$$

Put another way, the learned parameters  $\theta$  evolve approximately according to the linear dynamical system defined by the iteration matrix  $A_\Phi$ . As might be expected, TD(0) is stable if this linear dynamical system is globally stable in the usual sense (Borkar & Meyn, 2000).

The iteration matrix – and as we shall see, the global stability of the linear dynamical system – depends on the data distribution, the representation, and, to a lesser extent, on the discount factor. It does not, however, depend on the reward function, which only affects the accuracy of the TD fixed-point solution  $\theta_{TD}^*$ .

### 3.2. Stability Criteria

To understand the behaviour of TD(0), it is useful to contrast it with gradient descent on a weighted squared loss

$$\ell(\theta) = (\Phi\theta - \mathbf{y})^\top \Xi (\Phi\theta - \mathbf{y}),$$

where  $\mathbf{y}$  is a vector of supervised targets. Gradient descent on  $\ell(\theta)$  also corresponds to a linear dynamical system, albeit one whose iteration matrix is symmetric and positive definite. The behaviour of TD(0) is complicated by the fact that  $A_\Phi$  is not guaranteed to be positive definite or symmetric, as the matrix  $\Xi P^\pi$  itself is in general neither. In fact, the documented good behaviour of TD(0) arises in contexts where  $A_\Phi$  itself is closer to a gradient descent iteration matrix: positive definite when the data distribution is on-policy (Tsitsiklis & Roy, 1996), and symmetric when the Markov chain described by  $P^\pi$  is reversible (Ollivier, 2018).

Following a well-known result from linear system theory (see e.g. Zadeh & Desoer, 2008), the asymptotic behavior of TD(0) more generally depends on the eigenvalues of the iteration matrix.

**Proposition 3.1.** *TD(0) is stable if and only if the eigenvalues of the implied iteration matrix  $A_\Phi$  have positive real components, that is*

$$\text{Spec}(A_\Phi) \subset \mathbb{C}_+ := \{z : \text{Re}(z) > 0\}.$$

*We say that a particular choice of representation  $\Phi$  is **stable** for  $(P^\pi, \gamma, \Xi)$  when  $A_\Phi$  satisfies the above condition.*

*Proof.* See Appendix B for all proofs.  $\square$

Whenever the transition matrix, data distribution, and discount factor is evident, we will refer to  $\Phi$  simply as a stable representation.

### 3.3. Effect of Subspace Parametrization

When measuring the approximation error that arises from a particular representation  $\Phi$ , it suffices to consider the subspace spanned by the columns of  $\Phi$ . It therefore makes no difference whether these columns are orthogonal (corresponding, informally speaking, to correlated features) or not. By contrast, we now show that the stability of the learning process does depend on how the linear subspace spanned by  $\Phi$  is parametrized.

Recall that  $\Phi$  is orthogonal if  $\Phi^\top \Xi \Phi = I$ . As it turns out, the stability of an orthogonal representation is determined by the *induced transition matrix*  $\Pi P^\pi \Pi$ , which describes how next-state features affect the TD(0) value estimates.

**Proposition 3.2.** *An orthogonal representation  $\Phi$  is stable if and only if the real part of the eigenvalues of the induced transition matrix  $\Pi P^\pi \Pi$  is bounded above, according to*

$$\text{Spec}(\Pi P^\pi \Pi) \subset \{z \in \mathbb{C} : \text{Re}(z) < \frac{1}{\gamma}\}$$

*In particular,  $\Phi$  is stable if  $\rho(\Pi P^\pi \Pi) < \frac{1}{\gamma}$ .*

Although the original transition matrix satisfies the spectral radius condition with  $\rho(P^\pi) = 1$ , the induced transition matrix can have eigenvalues beyond the stable region and lead to learning instability.

More generally, a representation  $\Phi$  can be decomposed into an orthogonal basis and reparametrization  $\Phi = \Phi' R$ , where  $\Phi'$  is an orthogonal representation spanning the same space as  $\Phi$  and  $R \in \mathbb{R}^{d \times d}$  is a reparametrization for  $\Phi$ . The eigenvalues of the iteration matrix can be re-expressed as

$$\text{Spec}(A_\Phi) = \text{Spec}(R^\top A_{\Phi'} R) = \text{Spec}(R R^\top A_{\Phi'}).$$

Despite spanning the same space,  $\Phi$  and  $\Phi'$  have iteration matrices with different spectra:  $\text{Spec}(A_\Phi) \neq \text{Spec}(A_{\Phi'})$ . As a result, the stability of  $\Phi$  not only depends on the spectrum of  $A_{\Phi'}$ , but also how the reparametrization  $R$  shifts these eigenvalues. Put another way,  $\Phi$  may be unstable even if its orthogonal equivalent  $\Phi'$  is stable. The classical example of divergence given by Baird (1995) can be attributed to this phenomenon. In this example, the constructed representation expresses the same value functions as a stable tabular representation, but parametrizes the space in an different way and thus induces divergence.



### 3.4. Singular Vector Representations

The singular value decomposition is an appealing approach to representation learning: choosing vectors corresponding to large singular values guarantees, in a certain measure, low approximation error (Stachenfeld et al., 2014; Behzadian et al., 2019). Unfortunately, as now we show, doing so may be inimical to stability.

We denote  $\Phi_{SVD}$  and  $\Phi_{SR}$  the representations with the top  $d$  left singular vectors of  $P^\pi$  and  $\Psi$  as features. Recall that these vectors arise as part of a solution to a low-rank matrix approximation  $\|A - \hat{A}\|_\Xi$ . We write  $\hat{P}^\pi, \hat{\Psi} \in \mathbb{R}^{n \times n}$  to denote the corresponding rank- $d$  approximations.

**Proposition 3.3** (SVD). *The representation  $\Phi_{SVD}$  is stable if and only if the low-rank approximation  $\hat{P}^\pi$  satisfies*

$$\rho(\hat{P}^\pi) < \frac{1}{\gamma}.$$

**Proposition 3.4** (Successor Representation). *Recall that  $\text{Spec}(\Psi) \subset \mathbb{C}_+$ . The representation  $\Phi_{SR}$  is stable if and only if the low-rank approximation  $\hat{\Psi}$  satisfies*

$$\text{Spec}(\hat{\Psi}) \subset \mathbb{C}_+ \cup \{0\}.$$

Stability of a singular vector representation requires that the low-rank approximation maintain the spectral properties of the original matrix. This implies that such representations are *not stable* in general – the SVD low-rank approximation is chosen to minimize the norm of the error, and the spectrum of the approximation can deviate arbitrarily from the original matrix (Golub & van Loan, 2013). We note that the spectral conditions hold in the limit of almost-perfect approximation, but achieving this level of accuracy in practice may require an impractical number of additional features.

## 4. Representation Learning with Stability Guarantees

Our analysis of singular vector representations show that representations that optimize for alternative measures, such as approximation error, may lose properties of the transition matrix needed for stability. In this section, we study representations that are constrained, either in expressibility or in spectrum, to ensure stability.

### 4.1. Invariant Representations

We first consider representations whose induced transition matrix preserves the eigenvalues of the transition matrix to guarantee stability. These representations are closely linked to invariant subspaces of value functions that are closed under the transition dynamics of the policy.

**Definition 4.1.** *A representation  $\Phi$  is  $P^\pi$ -invariant if its corresponding linear subspace is closed under  $P^\pi$ , that is*

$$\text{Span}(P^\pi \Phi) \subseteq \text{Span}(\Phi).$$

$P^\pi$ -invariant subspaces are generated by the eigenspaces of  $P^\pi$ , and so invariant representations provide a natural way to reflect the geometry of the transition matrix. For these representations, we show that any eigenvalue of the induced transition matrix is also an eigenvalue of the transition matrix; this constraint ensures that invariant representations are always stable.

**Theorem 4.1.** *An orthogonal invariant representation  $\Phi$  satisfies*

$$\text{Spec}(\Pi P^\pi \Pi) \subseteq \text{Spec}(P^\pi) \cup \{0\}$$

*and is therefore stable.*

Parr et al. (2008) studied the quality of the TD fixed-point solution on invariant subspaces, and found it to directly correlate with how well the subspace models reward. Our findings on stability emphasize the importance of their result – with invariant representations that can predict reward, good value functions not only exist, but are also reliably discovered by TD(0).

Although estimation of eigenvectors for a nonsymmetric matrix is numerically unstable, finding orthogonal bases for their eigenspaces can be done tractably, for example through the Schur decomposition.

**Definition 4.2.** *Let  $A \in \mathbb{C}^{n \times n}$  be a complex matrix. A Schur decomposition of  $A$ , written  $\text{Schur}(A)$ , is  $URU^{-1}$ , where  $R$  is upper triangular and  $U = [u_1, u_2, \dots, u_n] \in \mathbb{C}^{n \times n}$  is orthogonal. For any  $k$ ,  $\text{Span}\{u_1, \dots, u_k\}$  is an  $A$ -invariant subspace.*

The Schur decomposition of  $P^\pi$  provides a sequence of vectors that span invariant subspaces, and can be constructed so that the first  $d$  basis vectors spans the top  $d$ -dimensional eigenspace of  $P^\pi$ . We define a representation using the first  $d$  Schur basis vectors to be the Schur representation.

When the transition matrix is reversible and data is on-policy, the Schur representation coincides with proto-value functions, and consequently also the successor representation (Machado et al., 2018). Unlike singular value representations, the Schur representation preserves the spectrum of the transition matrix at every step, and always guarantees stability.

**Corollary 4.1.1.** *The Schur representation is invariant and thus stable.*

A partial Schur basis can be constructed through orthogonal iteration, a generalized variant of power iteration.

**Proposition 4.1** (Golub & van Loan (2013)). *Let  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$  be the ordered eigenvalues of  $P^\pi$ . If  $|\lambda_d| > |\lambda_{d+1}|$  and  $\Phi_0 \in \mathbb{C}^{n \times d}$ , the sequence  $\Phi_1, \Phi_2, \dots$  generated via orthogonal iteration is*

$$\Phi_k = \text{ORTHOG}(\text{Span}(P^\pi \Phi_{k-1}))$$

where  $\text{ORTHO}(\cdot)$  finds an orthogonal basis. As  $k \rightarrow \infty$ ,  $\text{Span}(\Phi_k)$  converges to the unique top eigenspace of  $P^\pi$ .

In Section 5, we will see that the orthogonal iteration scheme can be approximated using a loss function and a target network (Mnih et al., 2015), and subsequently minimized with stochastic gradient descent, making it a potentially important tool for learning stable representations in practice.

## 4.2. Approximately Invariant Representations

In the previous section, we studied invariant representations, which are constrained to exactly preserve the eigenvalues of the transition matrix. We relax the notion of invariance to discuss a relaxation to approximate invariance, for which the spectrum of the induced matrix deviates from the transition matrix by a controlled amount, while still preserving stability. We find that approximate invariance leads to interesting implications for representations that span a Krylov subspace generated by rewards (Petrik, 2007; Parr et al., 2007).

**Definition 4.3.** A representation is  $\epsilon$ -invariant if

$$\max_{v \in \text{Span}(\Phi)} \frac{\|\Pi P^\pi v - P^\pi v\|_\Xi}{\|v\|_\Xi} \leq \epsilon.$$

An approximately invariant representation spans a space in which the transition dynamics are not fully closed, but approximately so, as measured by the  $\Xi$ -norm. We provide a simple condition of when an  $\epsilon$ -invariant representation is stable under assumptions of diagonalizability of the transition matrix. If  $P^\pi$  is diagonalizable with eigenbasis  $A \in \mathbb{C}^{n \times n}$ , the distance between the eigenvalues of the induced transition matrix  $\Pi P^\pi \Pi$  and the original transition matrix  $P^\pi$  can be bounded by a function of a)  $\epsilon$ , the degree of approximate invariance and b) the condition number of the eigenbasis  $\kappa_\Xi(A) = \|A\|_\Xi \|A^{-1}\|_\Xi$  (Trefethen & Embree, 2005).

**Theorem 4.2.** Let  $\Phi$  be an orthogonal and  $\epsilon$ -invariant representation for  $(P^\pi, \gamma, \Xi)$ . If  $P^\pi$  is diagonalizable with eigenbasis  $A$ , then  $\Phi$  is stable if

$$\epsilon < \frac{1 - \gamma}{\gamma} \frac{1}{\kappa_\Xi(A)}.$$

This bound is quite stringent, especially for discount factors close to one and ill-conditioned eigenvector bases, but may be improved if the transition matrix has a special structure. For the general setting when the transition matrix is not diagonalizable, similar but more complicated bounds exist (Shi & Wei, 2012).

Approximately invariant representations are of particular interest when studying the Krylov subspace generated by rewards,  $\mathcal{K}_d(P^\pi, r)$ .

$$\mathcal{K}_d(P^\pi, r) = \text{Span}\{r, P^\pi r, \dots, (P^\pi)^{d-1} r\}.$$

Representations that span this space admit a simple form of approximate invariance.

**Proposition 4.2.** A representation spanning  $\mathcal{K}_d(P^\pi, r)$  is  $\epsilon$ -invariant if

$$\frac{\|\Pi P^\pi v - P^\pi v\|_\Xi}{\|v\|_\Xi} \leq \epsilon$$

Where  $v = (I - \Pi_{d-1})(P^\pi)^{d-1} r$ , and  $\Pi_{d-1}$  is a projection onto the  $(d-1)$ -dimensional Krylov subspace  $\mathcal{K}_{d-1}(P^\pi, r)$ .

Orthogonal representations for this Krylov subspace are approximately invariant if they can predict the reward at the  $d+1$ -th timestep well from the rewards attained in the first  $d$  timesteps. For rewards that diffuse through the environment rapidly and can be predicted easily, an orthogonal basis of the Krylov space generated by rewards is approximately invariant and thus stable. Challenging environments with sparse rewards and temporal separation however may require a prohibitively large Krylov space to guarantee stability. Note that there is an important distinction between orthogonal representations spanning a Krylov subspace and the Krylov basis itself: for most practical applications, rewards are highly correlated and because of the challenges of parametrization, the latter can be unstable.

## 4.3. Positive-Definite Representations

Invariant representations are stable because the spectrum of the projected transitions is constrained to closely mimic the eigenvalues of the transition matrix. What we call *positive definite representations* instead guarantee stability by constraining the set of expressible value functions to lie within a safe set. Positive definite representations are stable regardless of parametrization, unlike any family of representations discussed so far.

**Definition 4.4.** The set of positive-definite value functions  $\mathcal{S}_{PD} \subset \mathbb{R}^n$  is

$$\mathcal{S}_{PD} = \{v \in \mathbb{R}^n \mid \langle v, P^\pi v \rangle_\Xi < \gamma^{-1} \|v\|_\Xi^2\}.$$

Note that  $\mathcal{S}_{PD}$  is not necessarily closed under addition. The two-state MDP presented by Tsitsiklis & Roy (1996) where TD(0) diverges can be interpreted through the lens of this set. For this example, the state representation only expresses value functions outside of  $\mathcal{S}_{PD}$ , which “grow” faster than  $\gamma^{-1}$ , and consequently leads to divergence. We focus on representations whose span falls within this set of safe value functions.

**Definition 4.5.** We say that a representation is *positive-definite* if

$$\text{Span}(\Phi) \subseteq \mathcal{S}_{PD}.$$

Note that a positive definite representation remains so under reparametrization, unlike the general case. In the special

case of on-policy data,  $\mathcal{S}_{PD} = \mathbb{R}^n$  and all representations are positive-definite (Tsitsiklis & Roy, 1996).

**Theorem 4.3.** *A positive-definite representation  $\Phi$  has a positive-definite iteration matrix  $A_\Phi$ , and is thus stable.*

The Laplacian representation, which computes the spectral eigendecomposition of the symmetrized transition matrix

$$K := \frac{1}{2} \left( P^\pi + \Xi^{-1} P^\pi \Xi \right) = U \Lambda U^\top \Xi,$$

provides an interesting bifurcation of value functions into those that are positive-definite and those that are not. As a consequence of Theorem 4.3, a stable representation is obtained by using eigenvectors corresponding to eigenvalues smaller or equal to  $\frac{1}{\gamma}$ .

**Proposition 4.3.** *Let  $\lambda_1, \dots, \lambda_n$  be the eigenvalues of  $K$ , in decreasing order, and  $u_1, \dots, u_n$  the corresponding eigenvectors. Define  $d^*$  as the smallest integer such that  $\lambda_{d^*} < \frac{1}{\gamma}$ . For any  $i \leq n - d^*$ , the safe Laplacian representation  $\Phi$ , defined as*

$$\Phi = [u_{d^*}, u_{d^*+1}, \dots, u_{d^*+i}],$$

*is positive-definite and stable.*

While including eigenvectors for larger eigenvalues does not guarantee divergence, the basis  $[u_1, \dots, u_i]$  for  $i < d^*$  is unstable (See appendix). When the data is on-policy, all eigenvalues of  $K$  are below the threshold  $\frac{1}{\gamma}$ , and the safe Laplacian corresponds exactly to the original representation.

We finish our discussion with a cautionary point. Although positive-definite representations admit amenable optimization properties, such as invariance to reparametrization and monotonic convergence, they can only express value functions that satisfy a growth condition. Under on-policy sampling this growth condition is nonrestrictive, but as the policy deviates from the data distribution, the expressiveness of positive-definite representations reduces greatly.

## 5. Experiments

We complement our theoretical results with an experimental evaluation, focusing on the following questions:

- How closely do the theoretical conditions we describe match stability requirements in practice?
- Can stable representations be learned using samples?
- Can they be learned using neural networks?

We conduct our study in the four-room domain (Sutton et al., 1999). We augment this domain with a task where the agent must reach the top right corner to receive a reward of +1 (Figure 1). The policy evaluation problem is to accurately

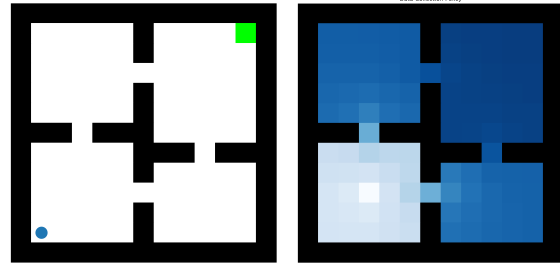


Figure 1. Left: The four room domain. Right: Data distribution

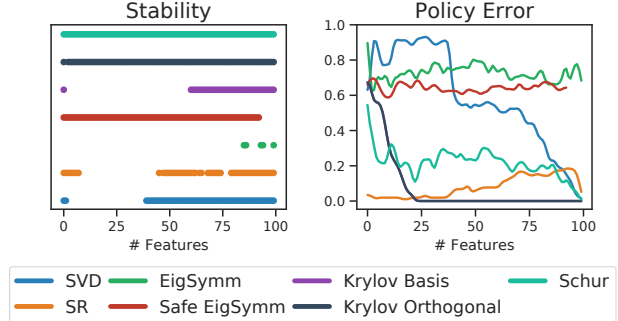


Figure 2. Stability (left) and approximation error (right) for different representation learning objectives. For stability, a marker is placed at  $d$  if the first  $d$  basis vectors forms a stable representation.

estimate the value function of a near-optimal policy from data consisting of trajectories sampled by an uniform policy.

We are interested in the usefulness of the representation learning schemes summarized in Table 1 as a function of the number of features  $d$  that are used. We measure both the stability of the learned representation and its accuracy in estimating the greedy policy with respect to the fixed value function. We chose the latter measure as it is more informative than value approximation error when the number of features is small. See Appendix C for full details about the experimental setup.

**Exact Representations:** We first consider the quality of the representations in exact form, assuming access to the true transition matrix and reward function (Figure 2). We find that the general empirical profiles for stability match our theoretical characterizations. Singular vectors of the successor representation have low error but are unstable for most choices of small  $d$ . Although the Krylov basis of rewards and its orthogonalization both have the same estimation errors, they have drastically different stability profiles, confirming our analysis from Section 4.2. Amongst the proposed methods that consistently produce stable representations, the Schur basis admits low error and with enough features, is fully expressible. In contrast, the safe Laplacian representation takes an irrecoverable performance hit, as it discards the top eigenvectors of the symmetrized transition matrix that contain reward-relevant information.

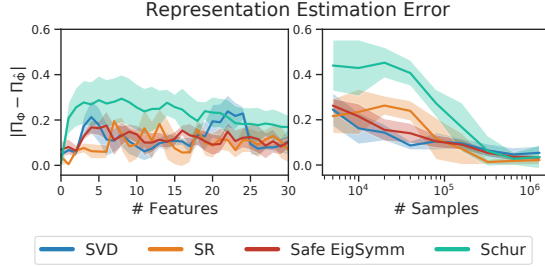


Figure 3. Learnability. We measure the difference between the true representation and one learned from an empirical transition matrix constructed from samples. Left: Error with 50000 transitions varying the number of features. Right: Error as the number of transitions varies when learning the first 10 features.

**Estimation with Samples:** In practice, representations must be learned from finite data. To test the numerical robustness of the representation learning schemes, we construct an empirical transition matrix from a variable number of samples and learn a representation using this matrix.

We measure the difference between the subspaces spanned by the estimated and true representation (Figure 3). We find that estimating the Schur representation can be more challenging than the other methods, and requires an order of magnitude more data to accurately compute than representations for singular vectors and spectral decompositions. This is a well-known problem in numerical linear algebra, as eigenspaces for nonsymmetric matrices (SCHUR) are more sensitive to perturbation and estimation error than for eigenspaces of symmetric matrices (SPECTRAL, SVD). This implies a three-way tradeoff between stability, approximation error, and ease of estimation when choosing a representation for a general environment. The successor representation is unstable, the safe Laplacian is limited in its approximation power, and the Schur decomposition is harder to learn from samples. The orthogonal Krylov basis emerges as a strong method by these measures, but requires additional knowledge in the guise of the reward function.

**Estimation with Neural Networks:** In our final set of experiments, we show that the Schur representation and the orthogonal Krylov representation can be learned by neural networks by performing stochastic gradient descent on certain auxiliary objectives.

It has been noted previously that training a representation network with a final linear layer to predict features causes the neural network to learn a basis for the target features (Bellemare et al., 2019). A  $d$ -dimensional Krylov representation then can be learned by predicting reward values at the next  $d$  time-steps. Similarly, orthogonal iteration for learning the Schur representation (Proposition 3.2) can be approximated with a two-timescale algorithm that (a) at each step, predicts the feature values of a fixed target representation network at the *next time step* and (b) infrequently refreshes the target representation network with the current.

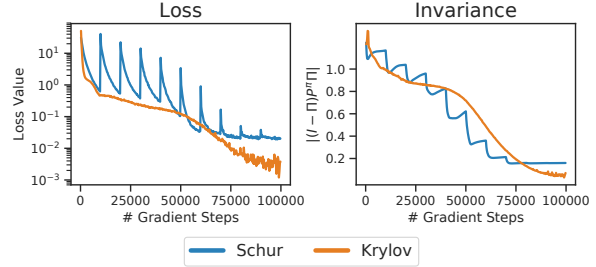


Figure 4. Learning to predict future rewards (Krylov) or future feature values (Schur) discovers approximately invariant stable representations.

As our stability guarantees hold for orthogonal representations, the neural network must learn uncorrelated features, which can be enforced explicitly or with a penalty-based orthogonality loss (Wu et al., 2018). We fully describe the auxiliary objectives and provide implementation details in Appendix C.

Figure 4 demonstrates that these predictive losses can be optimized easily with neural networks and can learn stable approximately invariant representations. We note that this auxiliary task of predicting future latent states has been heuristically proposed before (François-Lavet et al., 2018; Gelada et al., 2019), as a way to improve approximation errors. Our results indicate that such auxiliary tasks may not only help reduce approximation error, but more importantly, can mitigate divergence in the learning process and provide for stable optimization.

## 6. Conclusion

We have presented an analysis of stability guarantees for value-function learning under various representation learning procedures. Our analysis provides conditions for stability of many algorithms that learn features from transitions, and demonstrates how representation learning procedures constrained to respect the geometry of the transition matrix can induce stability. We demonstrated that the Schur decomposition and orthogonal Krylov bases are rich representations that mitigate divergence in off-policy value function learning, and further showed that they can be learned using stochastic gradient descent on a loss function.

Our work provides formal evidence that representation learning can prevent divergence without sacrificing approximation quality. To carry our results to the full practical case, stability should be extended to the sequence of policies that are encountered during policy iteration. One should also consider the effects of learning value functions and representations concurrently, and the ensuing interactions in the representation. Our work suggests that studying stable representations in these contexts can be a promising avenue forward for the development of principled auxiliary tasks for stable deep reinforcement learning.



## Acknowledgements

We thank Nicolas Le Roux, Marlos C. Machado, Courtney Paquette, Fabian Pedregosa, Doina Precup, and Ahmed Touati for helpful discussions and contributions. We additionally thank Marlos C. Machado and Courtney Paquette for constructive feedback on an earlier manuscript.

## References

- Baird, L. C. Residual algorithms: Reinforcement learning with function approximation. In *ICML*, 1995.
- Barnard, E. Temporal-difference methods and markov models. *IEEE Transactions on Systems, Man, and Cybernetics*, 1993.
- Behzadian, B., Gharatappeh, S., and Petrik, M. Fast feature selection for linear value function approximation. In *ICAPS*, 2019.
- Bellemare, M. G., Dabney, W., and Munos, R. A distributional perspective on reinforcement learning. In *ICML*, 2017.
- Bellemare, M. G., Dabney, W., Dadashi, R., Taïga, A. A., Castro, P. S., Roux, N. L., Schuurmans, D., Lattimore, T., and Lyle, C. A geometric perspective on optimal representations for reinforcement learning. In *Advances in Neural Information Processing Systems*, 2019.
- Benveniste, A., Priouret, P., and Métivier, M. *Adaptive Algorithms and Stochastic Approximations*. Springer-Verlag, Berlin, Heidelberg, 1990. ISBN 0387528946.
- Bertsekas, D. P. Approximate policy iteration: a survey and some new methods. *Journal of Control Theory and Applications*, 9:310–335, 2011.
- Bertsekas, D. P. Feature-based aggregation and deep reinforcement learning: a survey and some new implementations. *IEEE/CAA Journal of Automatica Sinica*, 6:1–31, 2018.
- Bodnar, C., Li, A., Hausman, K., Pastor, P., and Kalakrishnan, M. Quantile QT-opt for risk-aware vision-based robotic grasping. *arXiv*, 2019.
- Borkar, V. S. and Meyn, S. P. The o.d.e. method for convergence of stochastic approximation and reinforcement learning. *SIAM J. Control and Optimization*, 38:447–469, 2000.
- Bradtke, S. J. and Barto, A. G. Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22:33–57, 1996.
- Cabi, S., Colmenarejo, S. G., Novikov, A., Konyushkova, K., Reed, S., Jeong, R., Zolna, K., Ayta, Y., Budden, D., Vecerik, M., Sushkov, O., Barker, D., Scholz, J., Denil, M., de Freitas, N., and Wang, Z. A framework for data-driven robotics. *arXiv*, 2019.
- Chung, W., Nath, S., Joseph, A. G., and White, M. Two-timescale networks for nonlinear value function approximation. In *International Conference on Learning Representations*, 2019.
- Dalal, G., Szörényi, B., Thoppe, G., and Mannor, S. Finite sample analyses for td(0) with function approximation. In *AAAI*, 2017.
- Dann, C., Neumann, G., and Peters, J. Policy evaluation with temporal differences: a survey and comparison. *J. Mach. Learn. Res.*, 15:809–883, 2014.
- Dayan, P. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 5:613–624, 1993.
- François-Lavet, V., Bengio, Y., Precup, D., and Pineau, J. Combined reinforcement learning via abstract representations. *arXiv*, 2018.
- Gelada, C., Kumar, S., Buckman, J., Nachum, O., and Bellemare, M. G. DeepMDP: Learning continuous latent space models for representation learning. In *Proceedings of the International Conference on Machine Learning*, 2019.
- Golub, G. H. and van Loan, C. F. *Matrix Computations*. JHU Press, fourth edition, 2013. ISBN 1421407949 9781421407944. URL <http://www.cs.cornell.edu/cv/GVL4/golubandvanloan.htm>.
- Gordon, G. J. Stable function approximation in dynamic programming. In *ICML*, 1995.
- Jaderberg, M., Mnih, V., Czarnecki, W., Schaul, T., Leibo, J. Z., Silver, D., and Kavukcuoglu, K. Reinforcement learning with unsupervised auxiliary tasks. *ArXiv*, abs/1611.05397, 2016.
- Lagoudakis, M. G. and Parr, R. Least-squares policy iteration. *J. Mach. Learn. Res.*, 4:1107–1149, 2003.
- Levine, N., Zahavy, T., Mankowitz, D., Tamar, A., and Mannor, S. Shallow updates for deep reinforcement learning. In *Advances in Neural Information Processing Systems*, 2017.
- Machado, M. C., Rosenbaum, C., Guo, X., Liu, M., Tesauro, G., and Campbell, M. Eigenoption discovery through the deep successor representation. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=Bk8ZcAxR->.

- Maei, H. R., Szepesvari, C., Bhatnagar, S., Precup, D., Silver, D., and Sutton, R. S. Convergent temporal-difference learning with arbitrary smooth function approximation. In *NIPS*, 2009.
- Mahadevan, S. and Maggioni, M. Proto-value functions: A laplacian framework for learning representation and control in markov decision processes. *J. Mach. Learn. Res.*, 8:2169–2231, 2007.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 2015.
- Ollivier, Y. Approximate temporal difference learning is a gradient descent for reversible policies. *ArXiv*, abs/1805.00869, 2018.
- Parr, R., Painter-Wakefield, C., Li, L., and Littman, M. L. Analyzing feature generation for value-function approximation. In *ICML '07*, 2007.
- Parr, R., Li, L., Taylor, G., Painter-Wakefield, C., and Littman, M. L. An analysis of linear models, linear value-function approximation, and feature selection for reinforcement learning. In *ICML '08*, 2008.
- Pathak, D., Agrawal, P., Efros, A. A., and Darrell, T. Curiosity-driven exploration by self-supervised prediction. *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 488–489, 2017.
- Petrik, M. An analysis of laplacian methods for value function approximation in mdps. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence, IJCAI'07*, pp. 2574–2579, San Francisco, CA, USA, 2007. Morgan Kaufmann Publishers Inc.
- Puterman, M. L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., USA, 1st edition, 1994. ISBN 0471619779.
- Shi, X. and Wei, Y. A sharp version of bauer–fike’s theorem. *Journal of Computational and Applied Mathematics*, 236(13):3218 – 3227, 2012. ISSN 0377-0427. doi: <https://doi.org/10.1016/j.cam.2012.02.021>. URL <http://www.sciencedirect.com/science/article/pii/S0377042712000787>.
- Stachenfeld, K. L., Botvinick, M., and Gershman, S. J. Design principles of the hippocampal cognitive map. In *Advances in Neural Information Processing Systems*, 2014.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT Press, 2nd edition, 2018.
- Sutton, R. S., Precup, D., and Singh, S. P. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artif. Intell.*, 112:181–211, 1999.
- Trefethen, L. N. and Embree, M. Spectra and pseudospectra : the behavior of nonnormal matrices and operators. 2005.
- Tsitsiklis, J. N. and Roy, B. V. Analysis of temporal-difference learning with function approximation. In *NIPS*, 1996.
- van Hasselt, H., Doron, Y., Strub, F., Hessel, M., Sonnerat, N., and Modayil, J. Deep reinforcement learning and the deadly triad. *ArXiv*, abs/1812.02648, 2018.
- Vecerik, M., Sushkov, O., Barker, D., Rothörl, T., Hester, T., and Scholz, J. A practical approach to insertion with variable socket position using deep reinforcement learning. 2019.
- Wu, Y., Tucker, G., and Nachum, O. The laplacian in rl: Learning representations with efficient approximations. *ArXiv*, abs/1810.04586, 2018.
- Yu, H. and Bertsekas, D. P. Basis function adaptation methods for cost approximation in mdp. In *Proceedings of the IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning*, 2009.
- Zadeh, L. and Desoer, C. *Linear system theory: the state space approach*. Courier Dover Publications, 2008.

## A. Linear Algebra and Spectral Theory

### A.1. Inner Products

A positive-definite symmetric matrix  $D \in \mathbb{R}^{k \times k}$  induces an inner product  $\langle \cdot, \cdot \rangle_D$  and norm  $\|\cdot\|_D$  on  $\mathbb{R}^k$ . Specifically, the inner product is written as  $\langle v, w \rangle_D = v^\top D w$ , and the corresponding norm  $\|v\|_D^2 = \langle v, v \rangle_D = v^\top D v$ . This corresponds to a Hilbert space  $(\mathbb{R}^k, \langle \cdot, \cdot \rangle_D)$ . In our work, we equip  $\mathbb{R}^n$  (where  $n = |\mathcal{S} \times \mathcal{A}|$ ) with the inner-product induced by the data distribution  $\Xi$ . We also equip  $\mathbb{R}^d$  (the parameter space) with the usual Euclidean inner product.

Most definitions and constructions with the Euclidean inner product generalize to arbitrary Hilbert spaces, some which we describe on  $\mathbb{R}^n$ . Two vectors  $v, w \in \mathbb{R}^n$  are *orthogonal* if  $\langle v, w \rangle_\Xi = v^\top \Xi w = 0$ . A matrix  $A \in \mathbb{R}^{n \times d}$  is *orthogonal* if the columns have unit norm, and are orthogonal to one another:  $A^\top \Xi A = I$ . The generalization of transposes and symmetric matrices comes through the adjoint of a matrix  $A \in \mathbb{R}^{n \times n}$ , written as  $A^* = \Xi^{-1} A^\top \Xi$ . A matrix is self-adjoint if  $A = A^*$ , and for matrices that are not self-adjoint, the symmetric component is given as  $\bar{A} = \frac{1}{2}(A + A^*)$ . We refer to  $\|A\|$  as the matrix norm induced by the equivalent norm on vectors.

Matrix decompositions for a matrix  $A \in \mathbb{R}^{n \times n}$  can be re-visited with respect to this inner-product.

- **Spectral Decomposition:** If  $A$  is self-adjoint, it admits a decomposition  $A = U \Lambda U^\top \Xi$ , where  $U \in \mathbb{R}^{n \times n}$  is an orthogonal matrix whose columns are eigenvectors of  $A$  and  $\Lambda$  a diagonal matrix with the corresponding eigenvalues.
- **SVD:**  $A$  admits a decomposition  $A = U \Sigma V^\top \Xi$ , where  $U \in \mathbb{R}^{n \times n}$  is an orthogonal matrix whose columns are the left singular vectors of  $A$ ,  $V \in \mathbb{R}^{n \times n}$  is an orthogonal matrix whose columns are the right singular vectors of  $A$ , and  $\Lambda$  a diagonal matrix with the corresponding singular values. Letting  $U_d, V_d \in \mathbb{R}^{n \times d}$  correspond to the first  $d$  singular vectors and  $\Sigma_d \in \mathbb{R}^{d \times d}$  the diagonal matrix with the corresponding singular values, then the low-rank approximation  $\hat{A} = U_d \Sigma_d V_d^\top \Xi$  minimizes  $\|A - \hat{A}\|_\Xi$  amongst all rank  $d$  matrices.

### A.2. Eigenvalues

We define the eigenvalues of  $A \in \mathbb{C}^{k \times k}$  to be the roots of the characteristic polynomial  $p(t) = \det(A - tI)$ . Some eigenvalues may correspond to a multiple root – we refer to this multiplicity as the algebraic multiplicity. Every eigenvalue  $\lambda$  corresponds to an eigenspace  $\mathcal{V}_\lambda$  of eigenvectors with this eigenvalue. If the algebraic multiplicity of any eigenvalue  $\lambda$  does not equal the dimensionality of  $\mathcal{V}_\lambda$ , then  $A$  is said to be *defective*. Otherwise, the matrix  $A$  is diagonalizable as  $PDP^{-1}$ , where  $P$  is a basis of eigenvectors of  $A$ , and  $D$  the corresponding eigenvalues.

We write  $\text{Spec}(A) = \{\lambda_1, \dots, \lambda_k\} \subset \mathbb{C}$  to denote the set of eigenvalues of the matrix  $A$ . The spectral radius of a matrix is the maximum magnitude of eigenvalues, written as  $\rho(A) = \sup_{\lambda \in \text{Spec}(A)} |\lambda|$ . For two matrices  $A \in \mathbb{C}^{k \times m}$ ,  $B \in \mathbb{C}^{m \times k}$ , we have the following cyclicity:  $\text{Spec}(AB) \setminus \{0\} = \text{Spec}(BA) \setminus \{0\}$ . As a consequence, we also have that  $\rho(AB) = \rho(BA)$ . We utilize this cyclicity heavily in the ensuing proofs.

The perturbation of eigenvalues for a diagonalizable matrix can be bounded simply via the Bauer-Fike theorem. Specifically, if  $A \in \mathbb{C}^{k \times k}$  is diagonalizable as  $PDP^{-1}$ , then eigenvalues of the perturbed matrix  $\lambda' \in \text{Spec}(A + E)$  can be bounded in distance from the original eigenvalues as  $\inf_{\lambda \in \text{Spec}(A)} |\lambda - \lambda'| \leq \|E\| \kappa(P)$ , where  $\kappa(P) = \|P\| \|P^{-1}\|$ . As a simple corollary of the Bauer-Fike Theorem, we have that  $\rho(A + E) \leq \rho(A) + \|E\| \kappa(P)$ .

## B. Proofs

**Proposition 3.1.** *TD(0) is stable if and only if the eigenvalues of the implied iteration matrix  $A_\Phi$  have positive real components, that is*

$$\text{Spec}(A_\Phi) \subset \mathbb{C}_+ := \{z : \text{Re}(z) > 0\}.$$

*We say that a particular choice of representation  $\Phi$  is **stable** for  $(P^\pi, \gamma, \Xi)$  when  $A_\Phi$  satisfies the above condition.*

*Proof of Proposition 3.1.* We review the update taken by TD(0) (equation 1), rewritten to express the connection to the implied iteration matrix  $A_\Phi = \Phi^\top \Xi (I - \gamma P^\pi) \Phi$ . Notice that  $A_\Phi \theta_{TD}^* = \Phi^\top \Xi r$ .

$$\begin{aligned} \theta_{k+1} - \theta_{TD}^* &= \theta_k - \eta (\Phi^\top \Xi (I - \gamma P^\pi) \Phi \theta_k - \Phi^\top \Xi r) - \theta_{TD}^* \\ &= \theta_k - \theta_{TD}^* - \eta (A_\Phi \theta_k - A_\Phi \theta_{TD}^*) \\ &= (I - \eta A_\Phi)(\theta_k - \theta_{TD}^*) \end{aligned}$$

Unrolling the iteration, the error to the optimal solution takes the form

$$\theta_k - \theta_{TD}^* = (I - \eta A_\Phi)^k (\theta_0 - \theta_{TD}^*)$$

This above iteration converges from any initialization  $\theta_0$  if and only if the spectral radius is bounded by one:  $\rho(I - \eta A_\Phi) < 1$ .

From here, we can easily show that TD(0) is stable if and only if  $\text{Spec}(A_\Phi) \subset \mathbb{C}_+$ . If there is some step-size  $\eta > 0$  for which  $\rho(I - \eta A_\Phi) < 1$ , then  $\text{Spec}(A_\Phi) \subset \mathbb{C}_+$ . Similarly, if  $\text{Spec}(A_\Phi) \subset \mathbb{C}_+$ , then letting  $\eta = \min_{\lambda \in \text{Spec}(A_\Phi)} \frac{\text{Re}(\lambda)}{|\lambda|^2}$  satisfies that  $\rho(I - \eta A_\Phi) < 1$ . □

**Proposition 3.2.** *An orthogonal representation  $\Phi$  is stable if and only if the real part of the eigenvalues of the induced transition matrix  $\Pi P^\pi \Pi$  is bounded above, according to*

$$\text{Spec}(\Pi P^\pi \Pi) \subset \{z \in \mathbb{C} : \text{Re}(z) < \frac{1}{\gamma}\}$$

*In particular,  $\Phi$  is stable if  $\rho(\Pi P^\pi \Pi) < \frac{1}{\gamma}$ .*

*Proof of Proposition 3.2.* For an orthogonal representation, the iteration matrix can be written as  $A_{TD}^\Phi = I - \gamma \Phi^\top \Xi P^\pi \Phi$ . Then,

$$\begin{aligned} \text{Spec}(A_\Phi) \subset \mathbb{C}_+ &\iff \text{Spec}(\Phi^\top \Xi P^\pi \Phi) \subset \{z \in \mathbb{C} : \text{Re}(z) < \frac{1}{\gamma}\} \\ &\iff \text{Spec}(\Pi P^\pi) \subset \{z \in \mathbb{C} : \text{Re}(z) < \frac{1}{\gamma}\} \\ &\iff \text{Spec}(\Pi P^\pi \Pi) \subset \{z \in \mathbb{C} : \text{Re}(z) < \frac{1}{\gamma}\} \end{aligned}$$

The second step falls from the cyclicity of the spectrum and the observation that for an orthogonal representation  $\Phi$ , the projection can be written as  $\Phi \Phi^\top \Xi = \Pi$ . The spectral radius condition is immediate. □

**Proposition 3.3 (SVD).** *The representation  $\Phi_{SVD}$  is stable if and only if the low-rank approximation  $\hat{P}^\pi$  satisfies*

$$\rho(\hat{P}^\pi) < \frac{1}{\gamma}.$$

*Proof of Proposition 3.3.* We can write the SVD factorization of the transition matrix as

$$P^\pi = \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{bmatrix} V_1^\top \\ V_2^\top \end{bmatrix} \Xi$$

Then, for  $\Phi_{SVD} = U_1$ ,  $\Pi P^\pi = U_1 \Sigma_1 V_1^\top \Xi = \hat{P}^\pi$ . The necessary and sufficient conditions follow from Proposition 3.2. □



**Proposition 3.4** (Successor Representation). *Recall that  $\text{Spec}(\Psi) \subset \mathbb{C}_+$ . The representation  $\Phi_{SR}$  is stable if and only if the low-rank approximation  $\hat{\Psi}$  satisfies*

$$\text{Spec}(\hat{\Psi}) \subset \mathbb{C}_+ \cup \{0\}.$$

*Proof of Proposition 3.4.* We can write the SVD factorization of the successor representation  $\Psi = (I - \gamma P^\pi)^{-1}$

$$\Psi = [U_1 \quad U_2] \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{bmatrix} \begin{bmatrix} V_1^\top \\ V_2^\top \end{bmatrix} \Xi \quad (I - \gamma P^\pi) = [V_1 \quad V_2] \begin{bmatrix} \Sigma_1^{-1} & 0 \\ 0 & \Sigma_2^{-1} \end{bmatrix} \begin{bmatrix} U_1^\top \\ U_2^\top \end{bmatrix} \Xi$$

Then, for  $\Phi_{SR} = U_1$ , the iteration matrix can be written as  $A_\Phi = U_1^\top \Xi V_1 \Sigma_1^{-1}$ .

Now, write  $\hat{\Psi}$  as  $U_1 \Sigma_1 V_1^\top \Xi$ , and denote  $\hat{\Psi}^+$  the Moore-Penrose pseudoinverse, written as  $\hat{\Psi}^+ = V_1 \Sigma_1^{-1} U_1^\top \Xi$ . Cyclicity of the spectrum shows that the eigenvalues of the iteration matrix  $A_\Phi$  coincide with those of  $\hat{\Psi}^+$ .

$$\text{Spec}(\hat{\Psi}^+) = \text{Spec}(V_1 \Sigma_1^{-1} U_1^\top \Xi) = \text{Spec}(U_1^\top \Xi V_1 \Sigma_1^{-1}) \cup \{0\} = \text{Spec}(A_\Phi) \cup \{0\}.$$

We obtain the result by recognizing that all the eigenvalues of  $\hat{\Psi}^+$  have positive real component iff the same is true for  $\hat{\Psi}$ :

$$\text{Spec}(\hat{\Psi}) \subset \mathbb{C}_+ \cup \{0\} \iff \text{Spec}(\hat{\Psi}^+) \subset \mathbb{C}_+ \cup \{0\}.$$

□

**Theorem 4.1.** *An orthogonal invariant representation  $\Phi$  satisfies*

$$\text{Spec}(\Pi P^\pi \Pi) \subseteq \text{Spec}(P^\pi) \cup \{0\}$$

*and is therefore stable.*

*Proof of Theorem 4.1.* Let  $\lambda$  be a nonzero eigenvalue of  $\Pi P^\pi \Pi$  with an eigenvector  $v$ . Since  $\Pi P^\pi \Pi v = \lambda v$ ,  $v \in \text{Span}(\Phi)$ .

Since  $P^\pi$  is invariant on  $\text{Span}(\Phi)$ ,  $P^\pi v = \lambda v$ , and therefore  $\lambda$  is an eigenvalue of  $P^\pi$ . Therefore,  $\text{Spec}(\Pi P^\pi \Pi) \subset \text{Spec}(P^\pi) \cup \{0\}$ .

The spectrum of  $P^\pi$  implies the stability of the representation.  $P^\pi$  is a stochastic matrix satisfying  $\rho(P^\pi) = 1$ , and thus  $\rho(\Pi P^\pi \Pi) \leq 1$ , implying stability through Proposition 3.2. □

**Proposition 4.1** (Golub & van Loan (2013)). *Let  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$  be the ordered eigenvalues of  $P^\pi$ . If  $|\lambda_d| > |\lambda_{d+1}|$  and  $\Phi_0 \in \mathbb{C}^{n \times d}$ , the sequence  $\Phi_1, \Phi_2, \dots$  generated via orthogonal iteration is*

$$\Phi_k = \text{ORTHOG}(\text{Span}(P^\pi \Phi_{k-1}))$$

*where  $\text{ORTHOG}(\cdot)$  finds an orthogonal basis. As  $k \rightarrow \infty$ ,  $\text{Span}(\Phi_k)$  converges to the unique top eigenspace of  $P^\pi$ .*

*Proof of Proposition 4.1.* See Theorem 7.3.1 in Golub & van Loan (2013). □

**Theorem 4.2.** *Let  $\Phi$  be an orthogonal and  $\epsilon$ -invariant representation for  $(P^\pi, \gamma, \Xi)$ . If  $P^\pi$  is diagonalizable with eigenbasis  $A$ , then  $\Phi$  is stable if*

$$\epsilon < \frac{1 - \gamma}{\gamma} \frac{1}{\kappa_\Xi(A)}.$$

*Proof of Theorem 4.2.* We can rewrite the definition of  $\epsilon$ -invariance in terms of a matrix norm:  $\|P^\pi \Pi - \Pi P^\pi \Pi\|_\Xi < \epsilon$ . Thus, letting  $E = \Pi P^\pi \Pi - P^\pi \Pi$ , we have  $\|E\|_\Xi < \epsilon$ .

Now, suppose that  $\Pi P^\pi \Pi$  has an eigenvalue, eigenvector pair  $(\lambda, v)$ . This means that  $v \in \text{Span}(\Phi)$ .

$$\lambda v = \Pi P^\pi \Pi v = P^\pi \Pi v + E v = P^\pi v + E v \implies \lambda \in \text{Spec}(P^\pi + E)$$

Now, the Bauer-Fike Theorem (see Appendix A above) thus implies that  $\rho(\Pi P^\pi \Pi) < \rho(P^\pi) + \epsilon \kappa_\Xi(A) < 1 + \epsilon \kappa_\Xi(A)$ . Now, if  $\epsilon < \frac{1 - \gamma}{\gamma} \frac{1}{\kappa_\Xi(A)}$ , then  $\rho(\Pi P^\pi \Pi) < \gamma^{-1}$ , and stability follows from Proposition 3.2. □

**Proposition 4.2.** A representation spanning  $\mathcal{K}_d(P^\pi, r)$  is  $\epsilon$ -invariant if

$$\frac{\|\Pi P^\pi v - P^\pi v\|_\Xi}{\|v\|_\Xi} \leq \epsilon$$

Where  $v = (I - \Pi_{d-1})(P^\pi)^{d-1}r$ , and  $\Pi_{d-1}$  is a projection onto the  $(d-1)$ -dimensional Krylov subspace  $\mathcal{K}_{d-1}(P^\pi, r)$ .

**Remark:** The vector  $v$  can be interpreted as the component of the reward at the  $d$ -th timestep that cannot be predicted from the first  $d-1$  timesteps.

*Proof of Proposition 4.2.* Any vector  $v \in \mathcal{K}_d(P^\pi, r)$  can be decomposed into two components:  $\Pi_{d-1}v + (I - \Pi_{d-1})v$ .

$$\begin{aligned} \frac{\|\Pi P^\pi v - P^\pi v\|_\Xi}{\|v\|_\Xi} &= \frac{\|\Pi P^\pi (\Pi_{d-1}v + (I - \Pi_{d-1})v) - P^\pi (\Pi_{d-1}v + (I - \Pi_{d-1})v)\|_\Xi}{\|\Pi_{d-1}v + (I - \Pi_{d-1})v\|_\Xi} \\ &= \frac{\|\Pi P^\pi (I - \Pi_{d-1}) - P^\pi (I - \Pi_{d-1})v\|_\Xi}{\|\Pi_{d-1}v\|_\Xi + \|(I - \Pi_{d-1})v\|_\Xi} \end{aligned}$$

This expression is maximized whenever  $v$  is nonzero and  $\|\Pi_{d-1}v\|_\Xi = 0$ , which is true whenever  $v = (I - \Pi_{d-1})(P^\pi)^{d-1}r$ .

$$\sup_{v \in \text{Span}(\Phi)} \frac{\|\Pi P^\pi v - P^\pi v\|_\Xi}{\|v\|_\Xi} = \frac{\|\Pi P^\pi v - P^\pi v\|_\Xi}{\|v\|_\Xi}$$

□

**Theorem 4.3.** A positive-definite representation  $\Phi$  has a positive-definite iteration matrix  $A_\Phi$ , and is thus stable.

*Proof of Theorem 4.3.* First, we show that the iteration matrix  $A_\Phi$  is positive-definite, and then show that this implies stability.

For any  $x \in \mathbb{R}^d$ , let  $v = \Phi x$ . Because  $\Phi$  is positive-definite,  $v \in \mathcal{S}_{PD}$ . Notice that rearranging the definition of positive definiteness implies that  $\langle v, (I - \gamma P^\pi)v \rangle_\Xi > 0$ .

$$x^\top A_{TD}^\Phi x = v^\top \Xi (I - \gamma P^\pi)v = \langle v, (I - \gamma P^\pi)v \rangle_\Xi > 0.$$

Now, we consider an eigenvalue  $\lambda$  of the iteration matrix  $A_\Phi$ , and a corresponding unit eigenvector  $x \in \mathbb{C}^d$ . Writing  $x = a + ib$  for  $a, b \in \mathbb{R}^d$ ,

$$\text{Re}(x^\top A_\Phi x) = \text{Re}((a - ib)^\top A_\Phi (a + ib)) = a^\top A_\Phi a + b^\top A_\Phi b > 0.$$

Noticing that  $\bar{x}^\top A_\Phi x = \lambda \bar{x}^\top x = \lambda$ , and therefore the real component of  $\lambda$  is positive,  $\text{Re}(\lambda) > 0$ .

□

**Proposition 4.3.** Let  $\lambda_1, \dots, \lambda_n$  be the eigenvalues of  $K$ , in decreasing order, and  $u_1, \dots, u_n$  the corresponding eigenvectors. Define  $d^*$  as the smallest integer such that  $\lambda_{d^*} < \frac{1}{\gamma}$ . For any  $i \leq n - d^*$ , the safe Laplacian representation  $\Phi$ , defined as

$$\Phi = [u_{d^*}, u_{d^*+1}, \dots, u_{d^*+i}],$$

is positive-definite and stable.

*Proof of Proposition 4.3.* We shall show that  $\text{Span}(\{u_{d^*}, u_{d^*+1}, \dots, u_n\}) \subseteq \mathcal{S}_{PD}$ , which implies the proposition.

$$\langle v, P^\pi v \rangle_\Xi = \langle v, \frac{1}{2}(P^\pi + \Xi^{-1}(P^\pi)^\top \Xi)v \rangle_\Xi$$

Consider some  $v \in \text{Span}(\{u_{d^*}, u_{d^*+1}, \dots, u_n\})$  which can be expressed as  $\sum_{k=d^*}^n \alpha_k u_k$ . We have

$$\begin{aligned}
 \langle v, P^\pi v \rangle_\Xi &= \langle v, \frac{1}{2}(P^\pi + \Xi^{-1}(P^\pi)^\top \Xi)v \rangle_\Xi \\
 &= \left\langle \sum_{k=d^*}^n \alpha_k u_k, \frac{1}{2}(P^\pi + \Xi^{-1}(P^\pi)^\top \Xi) \sum_{k=d^*}^n \alpha_k u_k \right\rangle_\Xi \\
 &= \left\langle \sum_{k=d^*}^n \alpha_k u_k, \sum_{k=d^*}^n \lambda_k \alpha_k u_k \right\rangle_\Xi \\
 &< \gamma^{-1} \left\langle \sum_{k=d^*}^n \alpha_k u_k, \sum_{k=d^*}^n \alpha_k u_k \right\rangle_\Xi \\
 &= \gamma^{-1} \|v\|_\Xi^2
 \end{aligned}$$

Hence,  $v \in \mathcal{S}_{PD}$  and  $\text{Span}(\{u_{d^*}, u_{d^*+1}, \dots, u_n\}) \subseteq \mathcal{S}_{PD}$ . The second-to-last line is a result of eigenvalues being bounded by  $\gamma^{-1}$ .

Since  $\text{Span}(\Phi) \subseteq \text{Span}(\{u_{d^*}, u_{d^*+1}, \dots, u_n\})$ , we also have  $\text{Span}(\Phi) \subseteq \mathcal{S}_{PD}$ , and stability ensues from Theorem 4.3.

As a sidenote, we can use this same sequence of steps to show that a representation using only the top eigenvectors of  $K$  is always *not stable*. Defining the representation  $\Phi = [u_1, u_2, \dots, u_{d^*-1}]$ , and following the same set of steps yields that  $\langle v, P^\pi v \rangle > \gamma^{-1} \|v\|_\Xi^2$  for any  $v \in \text{Span}(\Phi)$ . This implies that for this representation, the iteration matrix  $A_\Phi$  is negative-definite, and has *all* eigenvalues with negative real component, therefore not stable.  $\square$

## C. Empirical Evaluation

### C.1. Experimental Setup

**Four-room Domain:** The four-room domain (Sutton et al., 1999) has 104 discrete states arranged into four “rooms”. At any state, the agent can take one of four actions corresponding to cardinal directions; if a wall blocks movement in the selected direction, the agent remains in place.

**Policy Evaluation:** We augment this domain with a task where the agent must reach the top right corner of the environment. The corresponding reward function is sparse, with the agent receiving +1 reward when it is in the desired state, and zero otherwise. The policy evaluation problem is to find the value function of a near-optimal policy in the environment Epsilon-Greedy( $\pi^*$ ,  $\epsilon = 0.1$ ), which takes the optimal action with probability 0.9, and a randomly selected action otherwise. Data is collected by rolling out 50-step trajectories from the center of the bottom-left room with a uniform policy, which samples actions uniformly at random. The discount factor is  $\gamma = 0.99$ .

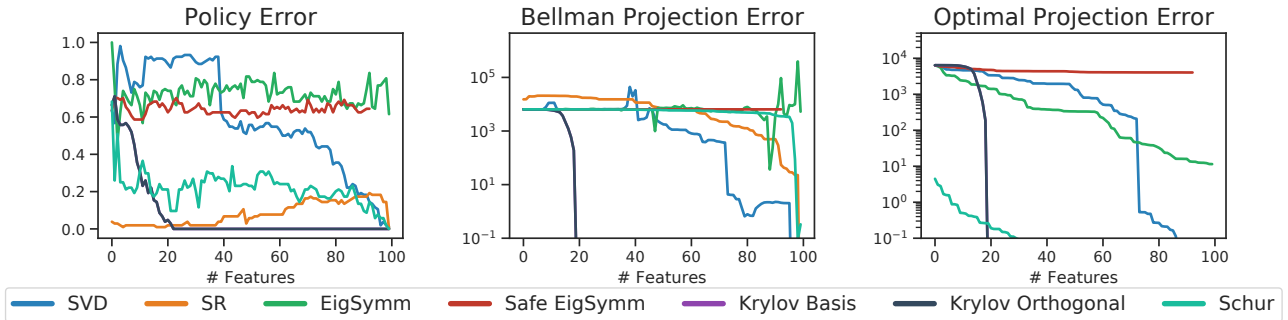
### C.2. Exact Evaluation

In this setting, the exact transition matrix  $P^\pi$  and data distribution  $\Xi$  are used to create the representation. We compute the decompositions according to Table 1 and Appendix A. Stability is measured for a given representation by explicitly creating the induced iteration matrix, computing the eigenvalues, and checking for real positive parts. To measure accuracy, we considered three metrics (Figure C.2).

- **Policy Accuracy: (displayed in paper)** This measures how well the greedy policy for the true value function matches the greedy policy for the estimated value function. This is given as

$$\frac{1}{|S|} \sum_{s \in S} \delta(\arg \max_a \hat{Q}(s, a) \neq \arg \max_a Q^\pi(s, a))$$

- **Optimal Projection Error:** This measures how far the true value function is from the subspace of expressible value functions  $\|Q^\pi - \Pi Q^\pi\|_\Xi$ . As the number of features increases, this error monotonically decreases, but may not be indicative of the quality of the solution.
- **Bellman Projection Error:** This measures how far the solution reached by TD(0) (the TD-fixed point) is from the true value function:  $\|Q^\pi - \Phi \theta_{TD}^*\|_\Xi$ . This measure of error is nonmonotonic (adding extra features can cause errors to increase) and unbounded. Furthermore, in the regime of a low number of features, this error greatly underestimates the quality of the recovered solution.



### C.3. Estimation from Samples

To measure how well the representations can be measured using samples, we consider the difference between the subspace spanned by the estimated and true representations. In particular, we sample  $t$  transitions from the data distribution, and reconstruct the empirical transition matrix  $\hat{P}^\pi$  given these transitions. If a particular  $(s, a)$  pair is never sampled, the prior we use for the transition matrix is that taking this action deterministically leads back to  $s$ . We construct the estimated



representation as  $\hat{\Phi}$ , and measure the distance between the true representation  $\Phi$  and the estimated representation  $\hat{\Phi}$  as  $\|\Pi_{\Phi} - \Pi_{\hat{\Phi}}\|_{\Xi, F}$ . The Frobenius norm  $\|\cdot\|_{\Xi, F}$  is selected in particular as this measures an expected distance, as compared to the maximum distance, measured by the operator norm  $\|\cdot\|_{\Xi}$ .

#### C.4. Estimation with Gradient Descent:

When learning the representation using gradient descent, we train a network  $f(s, a; \theta)$  with one hidden layer with  $d$  units with no activation function, that takes in state-action pairs encoded in one-hot form (as vectors in  $\mathbb{R}^{|\mathcal{S} \times \mathcal{A}|}$ ) and outputs in  $\mathbb{R}^d$ . In our experiments,  $d = 21$ . The value of the units in the hidden layer is the representation  $\phi(s, a; \theta)$ . The network is trained with a minibatch size of 32 for 100,000 steps, all implemented in Jax.

- **Schur Decomposition:** To mimic the orthogonal iteration procedure, we use the following training loss function, where  $\theta_t$  are the parameters for the target network.

$$\mathcal{L}(\theta; \theta_t) = \mathbb{E}_{\substack{(s,a) \sim \xi \\ s' \sim P(\cdot|s,a)}} \left[ \|f(s, a; \theta) - \mathbb{E}_{a' \sim \pi}[\phi(s', a'; \theta_t)]\|^2 \right]$$

This loss is optimized using stochastic gradient descent with a step-size of 4. The target network is updated every 10,000 steps, and after every target network update, the representation is renormalized to satisfy  $\mathbb{E}_{(s,a) \sim \xi} [\phi(s, a; \theta)_i^2] = 1$ .

- **Reward Krylov Basis:** We use the following regression training loss function

$$\mathcal{L}(\theta) = \mathbb{E}_{(s_1, a_1) \sim \xi} \left[ \sum_{i=1}^d \left( f(s, a; \theta)_i - \mathbb{E}_{(s_2, a_2, s_3, a_3, \dots, s_d, a_d) \sim P^{\pi}} [r(s_i, a_i)] \right)^2 \right]$$

where the inner expectation comes from trajectories that are generated from the policy  $\pi$  being evaluated starting from  $(s_1, a_1)$ . Although this loss requires that the evaluated policy be run in the environment, it serves a didactic purpose to show that these Krylov bases can be learned with additional domain knowledge. This loss is optimized using the Adam optimizer with a learning rate of  $10^{-3}$ .