
Reducing Sampling Error in Batch Temporal Difference Learning

Brahma S. Pavse¹ Ishan Durugkar¹ Josiah P. Hanna^{2,3} Peter Stone^{1,4}

Abstract

Temporal difference (TD) learning is one of the main foundations of modern reinforcement learning. This paper studies the use of TD(0), a canonical TD algorithm, to estimate the value function of a given policy from a batch of data. In this batch setting, we show that TD(0) may converge to an inaccurate value function because the update following an action is weighted according to the number of times that action occurred in the batch – not the true probability of the action under the given policy. To address this limitation, we introduce *policy sampling error corrected*-TD(0) (PSEC-TD(0)). PSEC-TD(0) first estimates the empirical distribution of actions in each state in the batch and then uses importance sampling to correct for the mismatch between the empirical weighting and the correct weighting for updates following each action. We refine the concept of a certainty-equivalence estimate and argue that PSEC-TD(0) is a more data efficient estimator than TD(0) for a fixed batch of data. Finally, we conduct an empirical evaluation of PSEC-TD(0) on three batch value function learning tasks, with a hyperparameter sensitivity analysis, and show that PSEC-TD(0) produces value function estimates with lower mean squared error than TD(0).

1. Introduction

Reinforcement learning (RL) (Sutton & Barto, 2018) algorithms have been applied to a variety of sequential-decision making problems such as robot manipulation (Kober et al., 2013; Gu et al., 2016) and autonomous driving (Sallab et al., 2017). Many RL algorithms learn an optimal control policy by estimating the *value function*, a function that gives the expected return from each state when following a particular

policy (Puterman & Shin, 1978; Bertsekas, 1987; Konda & Tsitsiklis, 2000). These algorithms require accurate value function estimation with finite data. A fundamental approach to value function learning is the temporal difference (TD) algorithm (Sutton, 1988).

In this work, we focus on improving the accuracy of the value function learned by *batch* TD, where TD updates for a value function are computed from a fixed batch of data. We show that batch TD(0) may converge to an inaccurate value function since it ignores the known action probabilities of the policy it is evaluating. For example, consider a single state in which the evaluation policy selects between action a_1 or a_2 with probability 0.5. If, in the finite batch of observed data, a_1 actually happens to occur twice as often as a_2 then TD updates following a_1 will receive twice as much weight as updates following a_2 , even though in expectation they should receive the same weight. We describe this finite-sample error in the value function estimate as *policy sampling error*. To correct for policy sampling error we propose to first estimate the maximum likelihood policy from the observed data and then use importance sampling (Precup et al., 2000a) to account for the mismatch between the frequency of sampled actions and their true probability under the evaluation policy. Variants of this technique have been successful in multi-armed bandits (Li et al., 2015; Narita et al., 2018; Xie et al., 2018), policy evaluation (Hanna et al., 2019), and policy gradient learning (Hanna & Stone, 2019). However, we are the first to show that this technique can be used to correct for policy sampling error in value function estimation and the first to show the benefit of importance sampling in on-policy value function estimation. We show that by using the available policy information, our approach is more data efficient than vanilla batch TD(0). We call our new value function learning algorithm *batch policy sampling error corrected-TD(0)* (PSEC-TD(0)).

The contributions of the paper are the following:

1. Show that the fixed point that batch TD(0) converges to for a given policy is inaccurate with respect to the true value function.
2. Introduce the batch PSEC-TD(0) algorithm that reduces the policy sampling error in batch TD(0).
3. Refine the concept of a certainty-equivalence estimate

¹The University of Texas at Austin ²School of Informatics, University of Edinburgh ³To be joining the Computer Sciences Department, University of Wisconsin–Madison ⁴Sony AI. Correspondence to: Brahma S. Pavse <brahmasp@cs.utexas.edu>.

for TD(0) (Sutton, 1988) and provide theoretical justification that batch PSEC-TD(0) is more data efficient than batch TD(0).

- Empirically analyze batch PSEC-TD(0) in the tabular and function approximation setting.

2. Background

This section introduces notation and formally specifies the batch value function learning problem.

2.1. Notation and Definitions

Following the standard MDPv1 notation (Thomas, 2015), we consider a Markov decision process (MDP) with state space \mathcal{S} , action space \mathcal{A} , reward function R , transition dynamics function P , and discount factor γ (Puterman, 2014). In any state s , an agent selects stochastic actions according to a policy π , $a \sim \pi(\cdot|s)$. After taking an action a in state s the agent transitions to a new state $s' \sim P(\cdot|s, a)$ and receives reward $R(s, a, s')$. We assume \mathcal{S} and \mathcal{A} to be finite; however, our experiments also consider infinite sized \mathcal{S} and \mathcal{A} . We consider the episodic, discounted, and finite horizon setting. The policy and MDP jointly induce a *Markov reward process* (MRP), in which the agent transitions between states s and s' with probability $P(s'|s)$ and receives reward $R(s, s')$. Finally, $\mathbf{x}(s) : \mathcal{S} \rightarrow \mathbb{R}^d$ gives a column feature vector for each state $s \in \mathcal{S}$.

We are concerned with computing the value function, $v^\pi : \mathcal{S} \rightarrow \mathbb{R}$, that gives the *value* of any state. The value of a particular state is the expected discounted return, i.e. the expected sum of discounted rewards when following policy π from that state:

$$v^\pi(s) := \mathbf{E}_\pi \left[\sum_{k=0}^L \gamma^k R_{t+k+1} \mid s_t = s \right], \forall s \in \mathcal{S} \quad (1)$$

where L is the terminal time-step and the expectation is taken over the distribution of future states, actions, and rewards under π and P .

2.2. Batch Value Prediction

This work investigates the problem of approximating v^{π_e} given a batch of data, \mathcal{D} , and an evaluation policy, π_e . Let a single episode, τ , be defined as $\tau := (s_0, a_0, r_0, s_1, \dots, s_{L_\tau-1}, a_{L_\tau-1}, r_{L_\tau-1})$, where L_τ is the length of the episode τ . The batch of data consists of m episodes, i.e., $\mathcal{D} := \{\tau_i\}_{i=0}^{m-1}$. The policy that generated the batch of data is called the *behavior policy*, π_b . If π_b is the same as π_e for all episodes then learning is said to be done *on-policy*; otherwise it is *off-policy*.

In batch value prediction, a value function learning algorithm uses a fixed batch of data to learn an estimate \hat{v}^{π_e} that approximates the true value function, v^{π_e} . In this work,

we introduce algorithmic and theoretical concepts with the linear approximation of v^{π_e} :

$$\hat{v}^{\pi_e}(s) := \mathbf{w}^T \mathbf{x}(s)$$

thus, in the linear case, we seek to find a weight vector \mathbf{w} , such that $\mathbf{w}^T \mathbf{x}(s)$ approximates the true value, $v^{\pi_e}(s)$. However, our empirical study also considers the non-linear approximation of v^{π_e} . The error of the predicted value function, \hat{v}^{π_e} , with respect to the true value function, v^{π_e} , is measured by calculating the mean squared value error between $v^{\pi_e}(s)$ and $\hat{v}^{\pi_e}(s) \forall s \in \mathcal{S}$ weighted by the proportion of time spent in each state under policy π_e , $d_{\pi_e}(s)$. Thus, we seek to find a weight vector \mathbf{w} that minimizes:

$$\text{MSVE}(\mathbf{w}) := \sum_{s \in \mathcal{S}} d_{\pi_e}(s) \left(v^{\pi_e}(s) - \mathbf{w}^T \mathbf{x}(s) \right)^2 \quad (2)$$

In this work, we compare data efficiency between two algorithms, X and Y , as follows:

Definition 1. *Data Efficiency.* A prediction algorithm X is more data efficient than algorithm Y if estimates from X have, on average, lower MSVE than estimates from Y for a given batch size.

2.3. Batch Linear TD(0)

A fundamental algorithm for value prediction is the single-step temporal difference learning algorithm, TD(0). Algorithm 1 gives pseudo-code for the batch linear TD(0) algorithm described by Sutton (1988).

Algorithm 1 Batch Linear TD(0) to estimate v^{π_e}

- Input: policy to evaluate π_e , behavior policy π_b , batch \mathcal{D} , linear value function, $\hat{v} : \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$, step-size $\alpha > 0$, convergence threshold $\Delta > 0$
 - Initialize: weight vector \mathbf{w}_0 arbitrarily (e.g.: $\mathbf{w}_0 := \mathbf{0}$), aggregation vector $\mathbf{u} := \mathbf{0}$, batch process counter, $i = 0$
 - while** $|\mathbf{w}_{i+1} - \mathbf{w}_i| \geq \mathbf{1} \cdot \Delta$ **do**
 - for** each episode, $\tau \in \mathcal{D}$ **do**
 - for** each transition, $(s, a, r, s') \in \tau$ **do**
 - $\hat{y} \leftarrow r + \gamma \mathbf{w}_i^T \mathbf{x}(s')$
 - $\rho \leftarrow \frac{\pi_e(a|s)}{\pi_b(a|s)}$ {for on-policy, $\pi_b = \pi_e$ }
 - $\mathbf{u} \leftarrow \mathbf{u} + [\rho \hat{y} - \mathbf{w}_i^T \mathbf{x}(s)] \mathbf{x}(s)$
 - end for**
 - end for**
 - $\mathbf{w}_{i+1} \leftarrow \mathbf{w}_i + \alpha \mathbf{u}$ {batch update}
 - $\mathbf{u} \leftarrow \mathbf{0}$ {clear aggregation}
 - $i \leftarrow i + 1$
 - end while**
-

Sutton (1988) proved that batch linear TD(0) converges to a fixed point in the on-policy case i.e. when $\pi_e = \pi_b$. An off-policy batch TD(0) algorithm uses importance sampling ratios to ensure that the expected update is the same as it would be if actions were taken with π_e instead of π_b (Precup

et al., 2000a). Unlike on-policy TD(0), off-policy TD(0) is *not* guaranteed to converge (Baird, 1995).

3. Convergence of Batch Linear TD(0)

In this section, we discuss the convergence of batch linear TD(0) to a fixed-point, the *certainty equivalence estimate* (CEE) for the underlying Markov reward process (MRP). We refine this concept to better reflect our objective of evaluating a policy in an MDP and then prove that batch TD(0) converges to an equivalent fixed point that ignores knowledge of the known evaluation policy, π_e , leading to inaccuracy in the value function estimate. This result motivates our proposed algorithm.

First, we introduce additional notation and assumptions. In this section, we assume that we are in the on-policy setting ($\pi_b = \pi_e$). Let $\hat{\mathcal{S}}$ be the set of states and $\hat{\mathcal{A}}$ be the set of actions that appear in \mathcal{D} and let $\bar{R}(s)$ be the mean reward received when transitioning from state s in the batch \mathcal{D} . Finally, if the notation includes a hat ($\hat{\cdot}$), it is the maximum-likelihood estimate (MLE) according to \mathcal{D} . For example, $\hat{\pi}$ is the MLE of π_b . Sutton (1988) proved that batch linear TD(0) converges to the CEE. That is, it converges to the exact value function of the maximum likelihood MRP according to the observed batch. This exact value function can be calculated using dynamic programming (Bellman, 2003; Bertsekas, 1987) with the MLE MRP transition function. We call this value function estimate the *Markov reward process certainty equivalence estimate* (MRP-CEE).

Definition 2. *Markov Reward Process Certainty Equivalence Estimate (MRP-CEE) Value Function.* The MRP-CEE is the value function \hat{v}_{MRP} that, $\forall s, s' \in \hat{\mathcal{S}}$, satisfies:

$$\hat{v}_{\text{MRP}}(s) = \bar{R}(s) + \gamma \sum_{k \in \hat{\mathcal{S}}} \hat{P}(s'|s) \hat{v}_{\text{MRP}}(s'). \quad (3)$$

Having now defined the MRP-CEE value function, we prove that batch TD(0) converges to the MRP-CEE value function. This fact was first proven by Sutton (1988) (see Theorem 3 of Sutton (1988)), however the original proof only considers rewards upon termination and no discounting. The extension to rewards per-step and discounting is straightforward, but to the best of our knowledge has not appeared in the literature before. Following Sutton’s proof (Sutton, 1988), we first prove the extension before extending the proof to an MDP, where the data inefficiency of TD(0) becomes clear. Proof details are in Appendix C.

Theorem 1 (Batch Linear TD(0) Convergence). *For any batch whose observation vectors $\{\mathbf{x}(s)|s \in \hat{\mathcal{S}}\}$ are linearly independent, there exists an $\epsilon > 0$ such that, for all positive $\alpha < \epsilon$ and for any initial weight vector, the predictions for linear TD(0) converge under repeated presentations of the batch with weight updates after each complete presentation to the fixed-point (3).*

In RL, the transitions of an MRP are a function of the behavior policy *and* transition dynamics distributions. That is $\forall s, s' \in \hat{\mathcal{S}}$:

$$\begin{aligned} \hat{P}(s'|s) &= \sum_{a \in \hat{\mathcal{A}}} \hat{\pi}(a|s) \hat{P}(s'|s, a), \\ \bar{R}(s) &= \sum_{a \in \hat{\mathcal{A}}} \hat{\pi}(a|s) \bar{R}(s, a) \end{aligned}$$

where $\bar{R}(s, a)$ is the mean reward observed in state s on taking action a . We define a new certainty-equivalence estimate that separates these two factors. We call this new value function estimate the *Markov decision process certainty equivalent estimate* (MDP-CEE).

Definition 3. *Markov Decision Process Certainty Equivalence Estimate (MDP-CEE) Value Function.* The MDP-CEE is the value function, \hat{v}_{MDP} , that, $\forall s, s' \in \hat{\mathcal{S}}$, satisfies:

$$\hat{v}_{\text{MDP}}(s) = \sum_{a \in \hat{\mathcal{A}}} \hat{\pi}(a|s) \left(\bar{R}(s, a) + \gamma \sum_{s' \in \hat{\mathcal{S}}} \hat{P}(s'|s, a) \hat{v}_{\text{MDP}}(s') \right) \quad (4)$$

Given the definitions of \hat{P} and \bar{R} , the MRP-CEE (Definition 2) and MDP-CEE (Definition 3) are equivalent. Theorem 2 gives the convergence of batch TD(0) to the MDP-CEE value function. Proof details are in Appendix D.

Theorem 2 (Batch Linear TD(0) Convergence). *For any batch whose observation vectors $\{\mathbf{x}(s)|s \in \hat{\mathcal{S}}\}$ are linearly independent, there exists an $\epsilon > 0$ such that, for all positive $\alpha < \epsilon$ and for any initial weight vector, the predictions for linear TD(0) converge under repeated presentations of the batch with weight updates after each complete presentation to the fixed-point (4).*

The MDP-CEE value function highlights two sources of estimation error in the value function estimate: $P \neq \hat{P}$ and/or $\pi_e \neq \hat{\pi}$. We describe the former as *transition sampling error* and the latter as *policy sampling error*. Transition sampling error may be unavoidable in a model-free setting since we do not know P . However, we do know π_e and can use this knowledge to potentially correct policy sampling error. In the next section, we present an algorithm that uses the knowledge of π_e to correct for policy sampling error and obtain a more accurate value function estimate.

4. Batch Linear PSEC-TD(0)

In this section, we introduce the batch policy sampling error corrected-TD(0) (PSEC-TD(0)) algorithm that corrects for the policy sampling error in batch TD learning. From Theorem 2, batch TD(0) converges to the value function for the maximum likelihood policy, $\hat{\pi}$, instead of π_e . Under this view, PSEC-TD(0) treats policy sampling error as an off-policy learning problem and uses importance sampling (Precup et al., 2000a) to correct the weighting of TD(0)

updates from $\hat{\pi}$ to π_e . Even though importance sampling is usually associated with off-policy learning, this approach is applicable in the on- and off-policy cases.

In addition to \mathcal{D} and π_e , we assume we are given a set of policies, Π . Batch PSEC-TD(0) first computes the maximum likelihood estimate of the behavior policy:

$$\hat{\pi} := \operatorname{argmax}_{\pi' \in \Pi} \sum_{\tau \in \mathcal{D}} \sum_{t=0}^{L_{\tau}-1} \log \pi'(a_t^\tau | s_t^\tau)$$

This estimation can be done in a number of ways. For example, in the tabular setting we could use the empirical count of actions in each state. This count-based approach is often intractable, and hence, in many problems of interest we must rely on function approximation. When using function approximation, the policy estimate can be obtained by minimizing a negative log-likelihood loss function. Once $\hat{\pi}$ is computed, the batch PSEC-TD algorithm is the same as Algorithm 1 with $\hat{\pi}$ replacing π_b in the importance sampling ratio. That is, for transition (s, a, r, s') in \mathcal{D} , the contribution to the weight update is $\mathbf{u} \leftarrow \mathbf{u} + [\hat{\rho}\hat{y} - \mathbf{w}_t^T \mathbf{x}(s)] \mathbf{x}(s)$, where $\hat{\rho} := \frac{\pi_e(a|s)}{\hat{\pi}(a|s)}$ is the PSEC weight (refer to Line 8 in Algorithm 1). Thus, PSEC makes an importance sampling correction from the *empirical* to the *evaluation* policy distribution.

4.1. Convergence of Batch Linear PSEC-TD(0)

Section 3 showed that batch TD(0) converges to two equivalent certainty-equivalence estimates. We now define a new certainty-equivalent estimate (CEE) to which our new batch PSEC-TD(0) algorithm converges. Intuitively, the MDP-CEE estimate (Definition 3) is the exact value function for the *MLE of the behavior policy*, $\hat{\pi}$, in the MLE of the MDP environment; our new algorithm converges to the exact value function for π_e in the MLE of the MDP environment, making it more data efficient than batch TD(0) once the batch size is large enough.

We define this new CEE as the *PSEC Markov Decision Process Certainty Equivalence Estimate* (PSEC-MDP-CEE) Value Function.

Definition 4. *PSEC Markov Decision Process Certainty Equivalence Estimate (PSEC-MDP-CEE) Value Function. The PSEC-MDP-CEE is the value function, $\hat{v}_{\text{PSEC-MDP}}^{\pi_e}$, that, $\forall s, s' \in \hat{\mathcal{S}}$, satisfies:*

$$\begin{aligned} \hat{v}_{\text{PSEC-MDP}}^{\pi_e}(s) &= \sum_{a \in \hat{\mathcal{A}}} \pi_e(a|s) [\hat{R}(s, a) \\ &+ \gamma \sum_{k \in \hat{\mathcal{S}}} \hat{P}(s'|s, a) \hat{v}_{\text{PSEC-MDP}}^{\pi_e}(s')] \end{aligned} \quad (5)$$

Theorem 3 states that batch PSEC-TD(0) converges to the new PSEC-MDP-CEE value function (Equation 5). Proof details are in Appendix E.

Theorem 3 (Batch Linear PSEC-TD(0) Convergence). *For*

any batch whose observation vectors $\{\mathbf{x}(s)|s \in \hat{\mathcal{S}}\}$ are linearly independent, there exists an $\epsilon > 0$ such that, for all positive $\alpha < \epsilon$ and for any initial weight vector, the predictions for linear PSEC-TD(0) converge under repeated presentations of the batch with weight updates after each complete presentation to the fixed-point (5).

We remark that convergence has only been shown for the on-policy setting. While PSEC-TD(0) can be applied in the off-policy setting, it may, like other semi-gradient TD methods, diverge when off-policy updates are made with function approximation (Baird, 1995). It is possible that combining PSEC-TD(0) with Emphatic TD (Mahmood et al., 2015) or Gradient-TD (Sutton et al., 2009) may result in provably convergent behavior with off-policy updates, however, that study is outside the scope of this work.

4.2. Extending PSEC to other TD Variants

In general, PSEC can improve any value function learning algorithm that computes the TD-error, δ , or equivalent errors. As an example, we consider the off-policy least-squares TD (LSTD) algorithm (Bradtke & Barto, 1996; Ghiassian et al., 2018), which analytically computes the exact parameters that minimize the TD-error in a batch of data using the following steps:

$$\begin{aligned} \mathbf{A} &= \sum_{(s,a,s') \in \mathcal{D}} [\hat{\rho} \mathbf{x}(s)(\mathbf{x}(s) - \gamma \mathbf{x}(s'))^T] \\ \mathbf{b} &= \sum_{(s,a,s') \in \mathcal{D}} R(s, a, s') \mathbf{x}(s) \\ \mathbf{w} &= \mathbf{A}^{-1} \mathbf{b}, \end{aligned}$$

where $\hat{\rho}$ is the PSEC weight. Even though we primarily consider TD(0) in this work, the extension to LSTD demonstrates that PSEC-TD can be extended to other value function learning algorithms.

5. Empirical Study

In this section, we empirically study PSEC-TD to answer the following questions:

1. Does batch PSEC-TD(0) lower MSVE compared to batch TD(0)?
2. Does batch linear PSEC-TD(0) empirically converge to its certainty-equivalence solution?
3. Does PSEC yield benefit when applied to LSTD?
4. What factors does PSEC's data efficiency depend on in the function approximation setting?

We briefly describe the RL domains used in our experiments.

- **Gridworld:** In this domain, an agent navigates a 4×4 grid to reach a corner. The *state and action spaces* are

discrete and we use a tabular representation for \hat{v}^{π_e} . PSEC-TD(0) uses count-based estimation for $\hat{\pi}$. The ground truth value function is computed with dynamic programming and the MSVE computation uniformly weights the error in each state. In Section 5.1, we consider a *deterministic* gridworld, where there is no transition dynamics sampling error.

- **CartPole:** In this domain, an agent controls a cart to balance a pole upright. The *state space is continuous and action space is discrete*. We only consider the on-policy setting. The evaluation policy is a neural network trained using REINFORCE (Williams, 1992). It has 2 hidden layers with 16 neurons. We evaluate PSEC with varying linear and neural network representations for the value function. $\hat{\pi}$ maps the raw state features to a softmax distribution over the actions with varying linear and neural network architectures. Since the true value function is unknown, we follow Pan et al. (2016) and use Monte Carlo rollouts from a fixed number of states sampled from episodes following the evaluation policy to approximate the ground-truth state-values of those states. We then compute the MSVE between the learned values and the average Monte Carlo return from these sampled states.
- **InvertedPendulum:** This domain is similar to CartPole, and the objective is the same – to balance a pole upright. However, the *state and action spaces are both continuous*. We only consider the on-policy setting. The evaluation policy is a neural network trained by PPO (Schulman et al., 2017). The network has 2 hidden layers with 64 neurons each. We evaluate PSEC with varying linear and neural network representations for the value function. The $\hat{\pi}$ estimate consists of two components: 1) a linear or neural network mapping from raw state features to the mean vector of a Gaussian distribution, and 2) parameters representing the log standard deviation of each element of the output vector. As in CartPole, we compute Monte Carlo rollouts for sampled states.

In all experiments, the value function learning algorithm iterates over the whole batch of data until convergence, after which the MSVE of the final value function is computed. Some experiments include a parameter sweep over the hyperparameters, which can be found in Appendix G.

5.1. Tabular Setting

In this set of experiments, we consider two variants of PSEC-TD that differ in the placement of the PSEC weight:

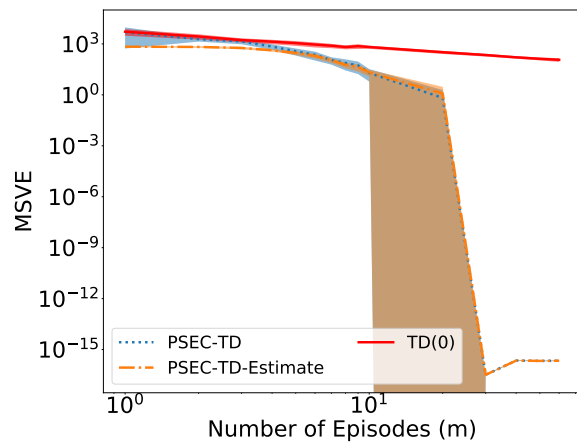
- **PSEC-TD-Estimate:** Multiplies $\hat{\rho}$ by the new estimate: $\hat{y} = R + \gamma w^T x(s')$.

- **PSEC-TD:** Multiplies $\hat{\rho}$ by the TD error: $\delta = (R + \gamma w^T x(s')) - w^T x(s)$.

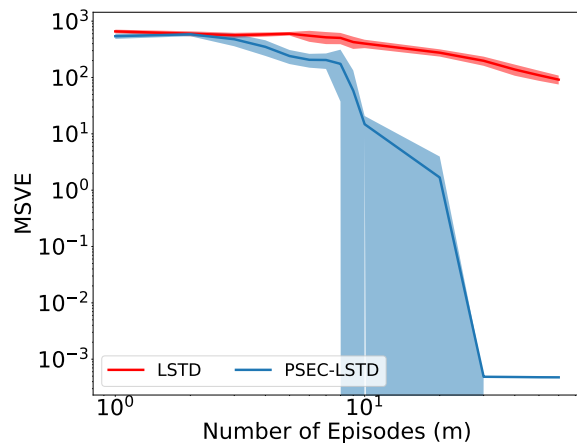
For off-policy TD(0), these placements of $\hat{\rho}$ are equivalent in expectation although the method using the TD-error has been reported to perform better in practice (Ghiassian et al., 2018). In this section, we focus on the on-policy results. Appendix F.1 includes off-policy results.

5.1.1. DATA EFFICIENCY

Figure 1 answers our first and third empirical questions, and shows that PSEC lowers MSVE compared to batch TD(0), and a variant of TD(0), LSTD(0). The gap between PSEC and its TD counterpart increases dramatically with more data; we discuss this observation in Section 5.1.2.



(a) On-policy (PSEC-)TD(0)

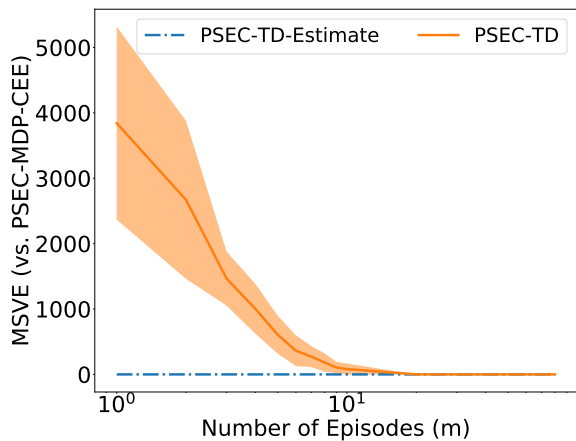


(b) On-Policy (PSEC-)LSTD(0)

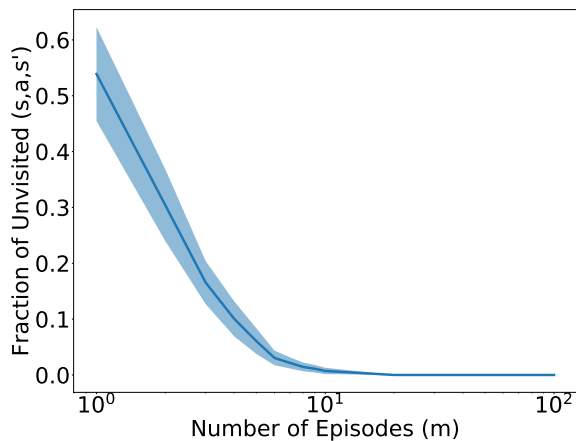
Figure 1. Deterministic Gridworld experiments. Both axes are log-scaled. Errors are computed over 200 trials with 95% confidence intervals. Asymmetric confidence intervals are due to log-scaling. Figure 1(a) and Figure 1(b) compare the data efficiency of PSEC-TD(0) and PSEC-LSTD(0) with their respective TD equivalents. Lower MSVE is better.

5.1.2. CONVERGENCE TO THE PSEC-CERTAINTY-EQUIVALENCE

To address our second empirical question, we empirically verify that both variants of batch linear PSEC, PSEC-TD and PSEC-TD-Estimate, converge to the dynamic programming computed PSEC-MDP-CEE value function (5) in Gridworld. According to Theorem 3, batch linear PSEC-TD-Estimate converges to the fixed-point (5) for all batch sizes.



(a)



(b)

Figure 2. Additional Gridworld experiments. Errors are computed over 50 trials with 95% confidence intervals. Figure 2(a) shows MSVE achieved by variants of linear batch PSEC-TD(0), PSEC-TD and PSEC-TD-Estimate, with respect to the PSEC-MDP-CEE (5). Figure 2(b) shows the fraction of unvisited (s, a, s') tuples.

We also empirically confirm that the other variant of PSEC, PSEC-TD converges to the same fixed-point (5) when the following condition holds true: only when all non-zero probability actions for each state in the batch have been sampled at least once. We note that when this condition is false, PSEC-TD-Estimate treats the value of taking that

action as 0. For example, if a state, s , appears in the batch and an action, a , that could take the agent to state s' does not appear in the batch, then PSEC-TD-Estimate treats the new estimate $R + \gamma \mathbf{w}^T \mathbf{x}(s')$ as 0, which is also done by the dynamic programming computation (5). We note that PSEC-TD converges to the fixed-point (5) only when this condition is true since the PSEC weight requires a fully supported probability distribution when applied to the TD-error estimate. From Figure 2(a) and Figure 2(b), we can see that this condition holds at batch size of 10 episodes. We also note that PSEC-TD(0) corrects policy sampling error for each (s, a, s') transition. Thus, when all such transitions are visited, PSEC fully corrects for all policy sampling error, which occurs at batch size of 10 episodes in this *deterministic* gridworld.

5.2. Function Approximation Setting

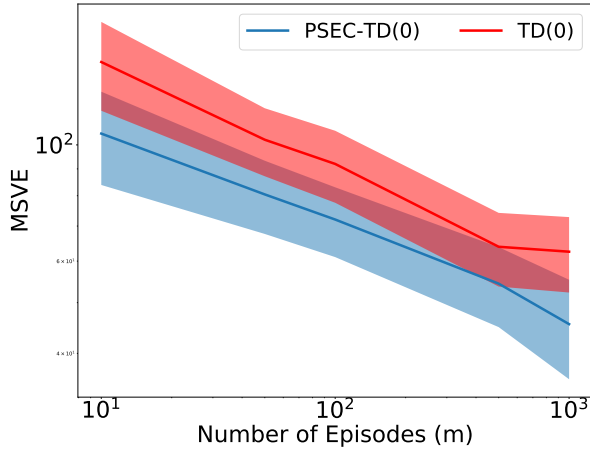
In this set of experiments, we answer our first and fourth empirical questions concerning function approximation in PSEC. Our experiments focus on applying only the second variant of PSEC, PSEC-TD, since we found that PSEC-TD-Estimate diverges. The results shown below are for the on-policy case. In addition to results of PSEC as a function of data size, we conduct experiments on a fixed batch size to better understand how components of the PSEC training process impact performance. Finally, we give a practical recommendation for use of batch PSEC-TD(0).

In these experiments, we have three function approximators: one for the value function; one to estimate the behavior policy; and the pre-learned behavior policy itself. When any are referred to as “fixed”, it means its architecture is unchanged. Due to space constraints, we only show a subset of results from CartPole and Inverted Pendulum; however, a fuller set of experiments can be found in Appendix F.2 and F.3. Note that in all PSEC training settings, PSEC performs gradient steps using the full batch of data, uses a separate batch of data as the validation data, and terminates training according to early stopping. Statistical significance is determined by Welch’s test (Welch, 1947) with a significance level of 0.05. For hyperparameter details refer to Appendix F.2 and F.3.

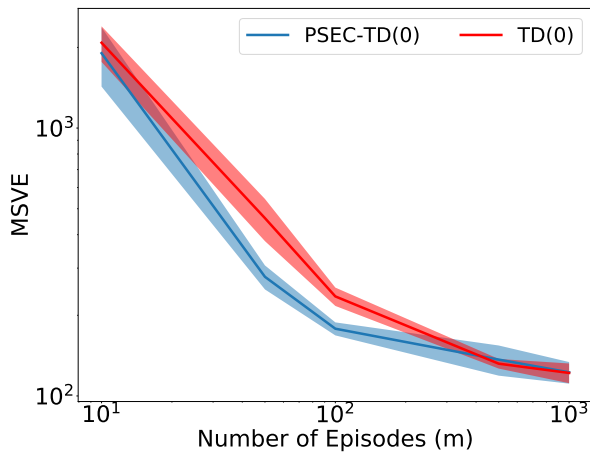
5.2.1. DATA EFFICIENCY

In CartPole, PSEC produced statistically significant improvement over TD in all batch sizes except 500. In InvertedPendulum, like in Gridworld, the improvement was marginal for smaller batch sizes, but produced statistically significant improvement with larger batch sizes. As data gets larger, we observe that both methods perform similarly for two reasons: 1) the PSEC weight approaches 1, which effectively becomes TD(0) and 2) saturation in value function representation capacity, which we discuss in Section

5.2.2. Note that while a thorough parameter sweep can achieve better performance, it is computationally expensive. The results shown here are with sweeps over only the value function model class and PSEC learning rate.



(a) CartPole

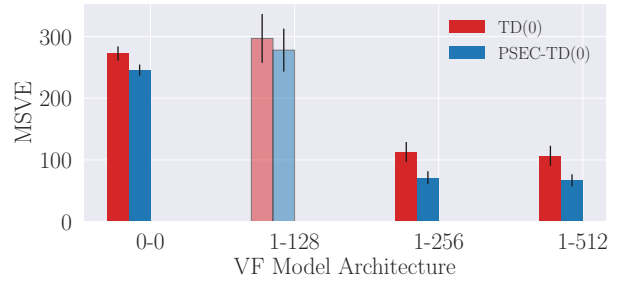


(b) InvertedPendulum

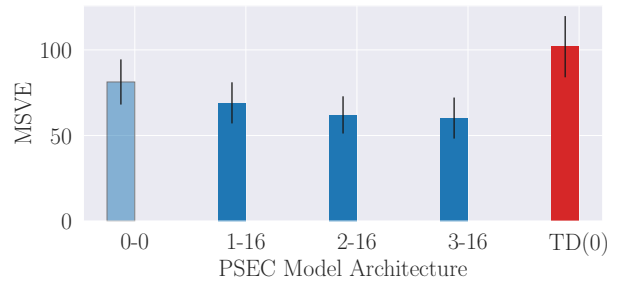
Figure 3. Comparing data efficiency of PSEC and TD on different batch sizes. Results for Figure 3(a) and Figure 3(b) are averaged over 400 and 250 trials resp. with shaded region of 95% confidence. Both axes are log-scaled. Lower MSVE is better.

5.2.2. ARCHITECTURE MODEL SELECTION

Figure 4(a) illustrates the impact of different value function classes on the data efficiency of TD and PSEC, while holding the PSEC model and behavior policy architectures fixed, on CartPole. We generally found that more expressive value function representations resulted in better data efficiency by both algorithms. We also found that the gap between PSEC and TD increased as the VF representation became more expressive. We hypothesize that even though PSEC finds a more accurate fixed point than TD in the space of



(a)



(b)

Figure 4. Figure 4(a) and Figure 4(b) compare data efficiency of PSEC, with varying VF model architectures, and PSEC, with varying model arch, respectively against TD on CartPole. Both use a batch size of 10 episodes, and results shown are averaged over 300 trials with error bars of 95% confidence. Darker shades represent statistically significant results. The label on the x axis shown is (# hidden layers - # neurons). Lower MSVE is better.

all value functions, the shown difference between the two algorithms is dependent on the space of representable value functions – a more representable function class can capture the difference between the two algorithms better. The lighter shades mean that any difference between PSEC and TD was statistically insignificant.

Figure 4(b) compares the data efficiency of PSEC against TD with varying PSEC neural network model architectures, while the value function and behavior policy architectures are fixed, on CartPole. In general, we found that more expressive network models produced better PSEC weights since they were able to better capture the MLE of the policy from the data. Unlike the NN PSEC policies, the linear function PSEC policy did not produce a statistically significant improvement over TD.

5.2.3. SENSITIVITY STUDIES

Due to space limitations, we defer the empirical analysis of other effects to Appendices F.2 and F.3. Figure 7 and Figure 12 indicate that a small learning rate for the PSEC model is preferred. Figure 9 and 14 indicate that some overfitting by

the PSEC model is tolerable, and perhaps, preferable, but extreme overfitting can degrade performance.

Practical Recommendation Based on our experiments, we recommend the following: 1) an expressive value function that can represent the more accurate fixed-point of PSEC-TD, 2) a PSEC model class that can represent the true behaviour policy but with awareness that extreme overfitting may hamper performance, and 3) a small learning rate.

6. Related Work

In this section, we discuss the literature on importance sampling with an estimated behavior policy and reducing sampling error in reinforcement learning.

The approach in this work has been motivated by prior work showing that importance sampling with an estimated behavior policy can lower variance when estimating an expected value in RL. [Hanna et al. \(2019\)](#) introduce a family of methods called regression importance sampling methods (RIS) and show that they have lower variance than importance sampling with the true behavior policy. [Hanna & Stone \(2019\)](#) show that a similar technique led to more sample-efficient policy gradient learning. These works are related to work in the multi-armed bandit ([Li et al., 2015](#); [Narita et al., 2018](#); [Xie et al., 2018](#)), causal inference ([Hirano et al., 2003](#); [Rosenbaum, 1987](#)), and Monte Carlo integration ([Henmi et al., 2007](#); [Delyon & Portier, 2016](#)) literature. In contrast, our work focuses on *value function learning*, where the focus is on learning the expected return at every state visited by the agent instead of across a set of actions (multi-armed bandit) or for some start states that are a subset of all the states the agent visits.

PSEC-TD(0) corrects policy sampling error through importance sampling with an estimated behavior policy. Other works avoid policy sampling error entirely by computing analytic expectations. Expected SARSA ([van Seijen et al., 2009](#)), learns action-values by analytically computing the expected return of the next state during bootstrapping as opposed to using the value of the sampled next action. The Tree-backup algorithm ([Precup et al., 2000b](#)) extends Expected SARSA to a multi-step algorithm. $Q(\sigma)$ ([Asis et al., 2017](#)) unifies SARSA ([Sutton, 1996](#); [Rummery & Niranjan, 1994](#)), Expected SARSA, and Tree-backups, to find a balance between sampling and analytic expectation computation. Our work is distinct from these in that we focus on learning state values which may be preferable for prediction as well as a variety of actor-critic approaches ([Konda & Tsitsiklis, 2000](#); [Mnih et al., 2016](#)). To the best of our knowledge, no other approach exists for correcting policy sampling error when learning state values.

7. Summary and Discussion

In batch value function approximation, we observed that TD(0) may converge to an inaccurate estimate of the value function due to policy sampling error. We proposed batch PSEC-TD(0) as a method to correct this error and showed that it leads to a more data efficient estimator than batch TD(0). In this paper, we theoretically analyzed PSEC-TD and empirically evaluated it in the tabular and function approximation settings. Our empirical study validated that PSEC converges to a more accurate fixed point than TD, and studied how the numerous components in the PSEC training setup impact its data efficiency with respect to TD.

Despite the data efficiency benefits that batch PSEC-TD(0) introduced, there are limitations. First, it requires knowledge of the evaluation policy, which on-policy TD(0) does not. This comparative disadvantage is only for the on-policy setting as both TD(0) and PSEC-TD(0) require knowledge of the evaluation policy for the off-policy setting. Additionally, PSEC-TD(0), in the off-policy case, has the advantage of not requiring knowledge of the behavior policy π_b . Second, the policy estimation step required by PSEC-TD(0) could potentially be computationally expensive. For instance, requiring the computation and storage of $\mathcal{O}(|S||A|)$ parameters in the tabular setting.

There are several directions for future work. First, our work focused on batch TD(0). We expect that a variant of PSEC can improve value function learning with n -step TD and TD(λ). Second, with an improved value function learning algorithm, it would be interesting to see if an agent can learn better control policies. Third, it would be interesting to theoretically and empirically study PSEC when learning the state-action values. Finally, automatically finding the optimal training setting for PSEC in the function approximation setting is another important direction for future work.

Acknowledgements

We thank Darshan Thaker, Elad Liebman, Reuth Mirsky, Sanmit Narvekar, Scott Niekum, Sid Desai, Yunshu Du, and the anonymous reviewers for reviewing our work and for their helpful comments. This work has taken place in the Learning Agents Research Group (LARG) at the Artificial Intelligence Laboratory, The University of Texas at Austin. LARG research is supported in part by grants from the National Science Foundation (CPS-1739964, IIS-1724157, NRI-1925082), the Office of Naval Research (N00014-18-2243), Future of Life Institute (RFP2-000), Army Research Office (W911NF-19-2-0333), DARPA, Lockheed Martin, General Motors, and Bosch. The views and conclusions contained in this document are those of the authors alone. Peter Stone serves as the Executive Director of Sony AI America and receives financial compensation for this work.

The terms of this arrangement have been reviewed and approved by the University of Texas at Austin in accordance with its policy on objectivity in research.

References

- Asis, K. D., Hernandez-Garcia, J. F., Holland, G. Z., and Sutton, R. S. Multi-step reinforcement learning: A unifying algorithm, 2017.
- Baird, L. Residual algorithms: Reinforcement learning with function approximation. In *Proceedings of the 12th International Conference on Machine Learning (ICML)*. 1995.
- Bellman, R. E. *Dynamic Programming*. Dover Publications, Inc., USA, 2003. ISBN 0486428095.
- Bertsekas, D. P. *Dynamic Programming: Deterministic and Stochastic Models*. Prentice-Hall, Inc., USA, 1987. ISBN 0132215810.
- Bradtke, S. and Barto, A. Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22: 33–57, 03 1996. doi: 10.1007/BF00114723.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym. *ArXiv*, abs/1606.01540, 2016.
- Delyon, B. and Portier, F. Integral approximation by kernel smoothing. *Bernoulli*, 22(4):2177–2208, Nov 2016. ISSN 1350-7265. doi: 10.3150/15-bej725. URL <http://dx.doi.org/10.3150/15-BEJ725>.
- Gerschgorin, S. Über die abgrenzung der eigenwerte einer matrix. *Izvestija Akademii Nauk SSSR, Serija Matematika*, 7(3):749–754, 1931.
- Ghiassian, S., Patterson, A., White, M., Sutton, R. S., and White, A. Online off-policy prediction, 2018.
- Glorot, X. and Bengio, Y. Understanding the difficulty of training deep feedforward neural networks. In Teh, Y. W. and Titterton, M. (eds.), *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pp. 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR. URL <http://proceedings.mlr.press/v9/glorot10a.html>.
- Gu, S., Holly, E., Lillicrap, T. P., and Levine, S. Deep reinforcement learning for robotic manipulation. *CoRR*, abs/1610.00633, 2016. URL <http://arxiv.org/abs/1610.00633>.
- Hanna, J. and Stone, P. Reducing sampling error in the monte carlo policy gradient estimator. In *Proceedings of the 18th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, May 2019.
- Hanna, J., Niekum, S., and Stone, P. Importance sampling policy evaluation with an estimated behavior policy. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, June 2019.
- Henmi, M., Yoshida, R., and Eguchi, S. Importance sampling via the estimated sampler. *Biometrika*, 94(4):985–991, 2007. URL <https://EconPapers.repec.org/RePEc:oup:biomet:v:94:y:2007:i:4:p:985-991>.
- Hirano, K., Imbens, G. W., and Ridder, G. Efficient estimation of average treatment effects using the estimated propensity score. *Econometrica*, 71(4):1161–1189, 2003. doi: 10.1111/1468-0262.00442. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/1468-0262.00442>.
- Kemeny, J. G., Snell, J. L., et al. *Finite markov chains*, volume 356. van Nostrand Princeton, NJ, 1960.
- Kingma, D. and Ba, J. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 12 2014.
- Kober, J., Bagnell, J., and Peters, J. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32:1238–1274, 09 2013. doi: 10.1177/0278364913495721.
- Konda, V. R. and Tsitsiklis, J. N. Actor-critic algorithms. In Solla, S. A., Leen, T. K., and Müller, K. (eds.), *Advances in Neural Information Processing Systems 12*, pp. 1008–1014. MIT Press, 2000. URL <http://papers.nips.cc/paper/1786-actor-critic-algorithms.pdf>.
- Li, L., Munos, R., and Szepesvári, C. Toward minimax off-policy value estimation. In *AISTATS*, 2015.
- Mahmood, A. R., Yu, H., White, M., and Sutton, R. S. Emphatic temporal-difference learning. *arXiv preprint arXiv:1507.01569*, 2015.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937, 2016.
- Narita, Y., Yasui, S., and Yata, K. Efficient counterfactual learning from bandit feedback. *CoRR*, abs/1809.03084, 2018. URL <http://arxiv.org/abs/1809.03084>.

- Pan, Y., White, A. M., and White, M. Accelerated gradient temporal difference learning. *CoRR*, abs/1611.09328, 2016. URL <http://arxiv.org/abs/1611.09328>.
- Precup, D., Sutton, R. S., and Singh, S. Eligibility traces for off-policy policy evaluation. In *Proceedings of the 17th International Conference on Machine Learning (ICML)*, pp. 759–766, 2000a.
- Precup, D., Sutton, R. S., and Singh, S. P. Eligibility traces for off-policy policy evaluation. In *International Conference on Machine Learning (ICML)*, 2000b.
- Puterman, M. L. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Puterman, M. L. and Shin, M. C. Modified policy iteration algorithms for discounted markov decision problems. *Management Science*, 24(11):1127–1137, 1978. ISSN 00251909, 15265501. URL <http://www.jstor.org/stable/2630487>.
- Rosenbaum, P. R. Model-based direct adjustment. *Journal of the American Statistical Association*, 82(398):387–394, 1987. ISSN 01621459. URL <http://www.jstor.org/stable/2289440>.
- Rummery, G. A. and Nanjund, M. On-line q-learning using connectionist systems. Technical report, 1994.
- Sallab, A. E., Abdou, M., Perot, E., and Yogamani, S. Deep reinforcement learning framework for autonomous driving. *Electronic Imaging*, 2017(19):70–76, 2017.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL <http://arxiv.org/abs/1707.06347>.
- Sutton, R. S. Learning to predict by the methods of temporal differences. *Machine learning*, 3(1):9–44, 1988.
- Sutton, R. S. Generalization in reinforcement learning: Successful examples using sparse coarse coding. In Touretzky, D. S., Mozer, M. C., and Hasselmo, M. E. (eds.), *Advances in Neural Information Processing Systems 8*, pp. 1038–1044. MIT Press, 1996.
- Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018. URL <http://incompleteideas.net/book/the-book-2nd.html>.
- Sutton, R. S., Maei, H. R., Precup, D., Bhatnagar, S., Silver, D., Szepesvári, C., and Wiewiora, E. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 993–1000, 2009.
- Thomas, P. S. A notation for markov decision processes. *CoRR*, abs/1512.09075, 2015. URL <http://arxiv.org/abs/1512.09075>.
- Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033, 2012.
- van Seijen, H., van Hasselt, H., Whiteson, S., and Wiering, M. A theoretical and empirical analysis of expected sarsa. In *2009 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning*, pp. 177–184, March 2009. doi: 10.1109/ADPRL.2009.4927542.
- Welch, B. L. The Generalization Of Student’s Problem When Several Different Population Variances Are Involved. *Biometrika*, 34(1-2):28–35, 01 1947. ISSN 0006-3444. doi: 10.1093/biomet/34.1-2.28. URL <https://doi.org/10.1093/biomet/34.1-2.28>.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Mach. Learn.*, 8(3?4):229?256, May 1992. ISSN 0885-6125. doi: 10.1007/BF00992696. URL <https://doi.org/10.1007/BF00992696>.
- Xie, Y., Liu, B., Liu, Q., Wang, Z., Zhou, Y., and Peng, J. Off-policy evaluation and learning from logged bandit feedback: Error reduction via surrogate policy. *CoRR*, abs/1808.00232, 2018. URL <http://arxiv.org/abs/1808.00232>.

A. Matrix Notation for Proofs

In this section, we introduce matrix-related notation for the proofs. Part of this notation is derived from Sutton (1988). We refer to state features using vectors indexed by the state. So features $x(i)$ for state i is referred to as x_i . Reward $r(j|i, a)$ (the reward for transitioning to state j from state i after taking action a) is referred to by r_{ij}^a . The policy $\pi(a|i)$ is π_i^a and the transition function $P(j|i, a)$ is accordingly the value p_{ij}^a . \mathcal{N} and \mathcal{T} are the set of non-terminal and terminal states respectively, and \mathcal{A} is the set of possible actions. \mathcal{D} is the batch of data used to train the value function. For any transition from a non-terminal state to a (non-) terminal state, $i \rightarrow j$, $1 \leq i \leq |\mathcal{N}|$ and $1 \leq j \leq |\mathcal{N} \cup \mathcal{T}|$. Finally, the maximum-likelihood estimate (MLE) of the above quantities according to \mathcal{D} , is given with a hat ($\widehat{\cdot}$) on top of the quantity. Further notations that are used in the proofs are explained when introduced.

Notation	Description	Dimension
\mathcal{S}	set of states	$ \mathcal{S} $
$\widehat{\mathcal{S}}$	set of states that appear in batch \mathcal{D}	$ \widehat{\mathcal{S}} $
\mathcal{A}	set of actions	$ \mathcal{A} $
$\widehat{\mathcal{A}}_i$	set of actions that appear in batch \mathcal{D} when agent is in state i	$ \widehat{\mathcal{A}}_i $
$\mathcal{N} \subset \mathcal{S}$	non-terminal states	$ \mathcal{N} $
$\widehat{\mathcal{N}} \subset \widehat{\mathcal{S}}$	non-terminal states that appear in batch \mathcal{D}	$ \widehat{\mathcal{N}} $
$\mathcal{T} \subset \mathcal{S}$	terminal states	$ \mathcal{T} $
$\widehat{\mathcal{T}} \subset \widehat{\mathcal{S}}$	terminal states that appear in batch \mathcal{D}	$ \widehat{\mathcal{T}} $
I	identity matrix	$ \mathcal{N} \times \mathcal{N} $
p_{jk}^π	probability of transitioning from state j to state k under a policy, π	-
p_{jk}^a	probability of transitioning from state j to state k after taking action a	-
\bar{r}_j	mean reward when transitioning from state j	-
\bar{r}_{ij}^a	mean reward when transitioning from state i to state j after taking action a	-
Q	$Q_{jk} := p_{jk}^\pi$	$ \mathcal{N} \times \mathcal{N} $
$[\mathbf{m}]_i$	expected reward on transitioning from state i to non-terminal state j i.e. $\sum_{j \in \mathcal{N}} p_{ij}^\pi r_{ij}$ or $\sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a r_{ij}^a$	$ \mathcal{N} $
$[\mathbf{h}]_i$	expected reward on transitioning from state i to non-terminal state j to terminal state i.e. $\sum_{j \in \mathcal{T}} p_{ij}^\pi r_{ij}$ or $\sum_{j \in \mathcal{T}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a r_{ij}^a$	$ \mathcal{N} $
$d_{\pi_e}(i)$	$\forall i \in \mathcal{S}$ weighted proportion of time spent in state i under policy π_e	-

B. Fixed-point for an MDP in the Per-step Reward and Discounted Case

In this section, we establish several fixed-points that we expect the value function to converge to in the discounted per-step reward case for an MRP and MDP. Note that Sutton (1988) derived these fixed-points for MRPs when rewards are only received on termination and there is no discounting. We first specify the fixed-points for an MRP in the discounted per-step reward case, and then extend this result for MDPs.

For both an MRP and MDP, we establish two types of fixed-points in the per-step reward and discounted case. The first fixed-point is the *true* fixed-point in that it is the value function computed assuming that we have access to the true policy and transition dynamics distributions. Ideally, we would like our value function learning algorithms to converge to this fixed-point. The second fixed-point is the *certainty-equivalence estimate* fixed-point, which is the value function computed using the maximum-likelihood estimates of the policy and transition dynamics from a batch of fixed data. We note that due to sampling error in the policy and transition dynamics, the certainty-equivalence estimate is an *inaccurate* estimate of the true value function. Finally, and only for MDPs, we specify the fixed-point that we expect the value function to learn after applying PSEC.

B.1. MRP True Fixed-Point

The true value function v for a state, $i \in \mathcal{N}$, induced by the policy-integrated transition dynamics p^π , reward function r , and policy, π , is given by:

$$\begin{aligned}
 v(i) &= \sum_{j \in \mathcal{N} \cup \mathcal{T}} p_{ij}^\pi [r_{ij} + \gamma v(j)] && \text{Bellman equation} \\
 &= \sum_{j \in \mathcal{N}} p_{ij}^\pi [r_{ij} + \gamma v(j)] + \sum_{j \in \mathcal{T}} p_{ij}^\pi r_{ij} && \text{expected return from } \mathcal{T}, v(\mathcal{T}) = 0 \\
 &= \sum_{j \in \mathcal{T}} p_{ij}^\pi r_{ij} + \sum_{j \in \mathcal{N}} p_{ij}^\pi \left[r_{ij} + \gamma \left[\sum_{k \in \mathcal{N} \cup \mathcal{T}} p_{jk}^\pi [r_{jk} + \gamma v(k)] \right] \right] && \text{recursively apply } v(i) \\
 v(i) &= \sum_{j \in \mathcal{T}} p_{ij}^\pi r_{ij} + \sum_{j \in \mathcal{N}} p_{ij}^\pi r_{ij} + \gamma \sum_{j \in \mathcal{N}} p_{ij}^\pi \sum_{k \in \mathcal{N} \cup \mathcal{T}} p_{jk}^\pi r_{jk} \\
 &\quad + \gamma^2 \sum_{j \in \mathcal{N}} p_{ij}^\pi \sum_{k \in \mathcal{N} \cup \mathcal{T}} p_{jk}^\pi v(k) \\
 &= \sum_{j \in \mathcal{T}} p_{ij}^\pi r_{ij} + \sum_{j \in \mathcal{N}} p_{ij}^\pi r_{ij} \\
 &\quad + \gamma \sum_{j \in \mathcal{N}} p_{ij}^\pi \sum_{k \in \mathcal{N}} p_{jk}^\pi r_{jk} + \gamma \sum_{j \in \mathcal{N}} p_{ij}^\pi \sum_{k \in \mathcal{T}} p_{jk}^\pi r_{jk} \\
 &\quad + \gamma^2 \sum_{j \in \mathcal{N}} p_{ij}^\pi \sum_{k \in \mathcal{N}} p_{jk}^\pi v(k) && \text{splitting } \mathcal{N} \text{ and } \mathcal{T}
 \end{aligned}$$

We define vectors, \mathbf{h} and \mathbf{m} with, $[\mathbf{h}]_i = \sum_{j \in \mathcal{T}} p_{ij}^\pi r_{ij}$, $[\mathbf{m}]_i = \sum_{j \in \mathcal{N}} p_{ij}^\pi r_{ij}$, and Q is the true transition matrix of the Markov reward process induced by π and P , i.e., $[Q]_{ij} = p_{ij}^\pi$. Then continuing from above, we have

$$v(i) = [\mathbf{h}]_i + [\mathbf{m}]_i + \gamma Q[\mathbf{h}]_i + \gamma Q[\mathbf{m}]_i + \gamma^2 Q^2[\mathbf{h}]_i + \gamma^2 Q^2[\mathbf{m}]_i + \dots \quad \text{unrolling } v(\mathcal{N} \cup \mathcal{T}) \quad (6)$$

$$= \left[\sum_{k=0}^{\infty} (\gamma Q)^k (\mathbf{m} + \mathbf{h}) \right]_i \quad (7)$$

$$v(i) = [(I - \gamma Q)^{-1}(\mathbf{m} + \mathbf{h})]_i \quad (8)$$

The existence of the limit and inverse are assured by Theorem A.1 in Sutton (1988). The theorem is applicable here since $\lim_{k \rightarrow \infty} (\gamma Q)^k = 0$.

B.2. MRP Certainty-Equivalence Fixed -Point

For the certainty-equivalence fixed-point, we consider a batch of data, \mathcal{D} . We follow the same steps and similar notation from Equation (8), with the slight modification that the maximum-likelihood estimate (MLE) of the above quantities according to \mathcal{D} , is given with a hat ($\hat{\cdot}$) on top of the quantity. The observed sets of non-terminal and terminal states in the batch are given by $\hat{\mathcal{N}}$ and $\hat{\mathcal{T}}$ respectively.

Then similar to above, we can derive the certainty-equivalence estimate of the value function according to the MLE of the MRP transition dynamics from the batch for a particular state i , $\forall i \in \hat{\mathcal{N}}$ is:

$$\hat{v}(i) = \left[(I - \gamma \hat{Q})^{-1}(\hat{\mathbf{m}} + \hat{\mathbf{h}}) \right]_i \quad (9)$$

B.3. MDP True Fixed-Point

The true value function, v^π , for a policy, π , for a state i , $\forall i \in \mathcal{N}$, induced by the transition dynamics and reward function, p and r is given by:

$$\begin{aligned}
 v^\pi(i) &= \sum_{a \in \mathcal{A}} \pi_i^a \sum_{j \in \mathcal{N} \cup \mathcal{T}} p_{ij}^a [r_{ij}^a + \gamma v^\pi(j)] && \text{Bellman equation} \\
 &= \sum_{a \in \mathcal{A}} \pi_i^a \sum_{j \in \mathcal{N}} p_{ij}^a [r_{ij}^a + \gamma v^\pi(j)] + \sum_{a \in \mathcal{A}} \pi_i^a \sum_{j \in \mathcal{T}} p_{ij}^a r_{ij}^a && \text{expected return from } \mathcal{T}, v^\pi(\mathcal{T}) = 0 \\
 \\
 v^\pi(i) &= \sum_{a \in \mathcal{A}} \pi_i^a \sum_{j \in \mathcal{T}} p_{ij}^a r_{ij}^a + \sum_{a \in \mathcal{A}} \pi_i^a \sum_{j \in \mathcal{N}} p_{ij}^a \left[r_{ij}^a + \gamma \left[\sum_{a' \in \mathcal{A}} \pi_j^{a'} \sum_{k \in \mathcal{N} \cup \mathcal{T}} p_{jk}^{a'} [r_{jk}^{a'} + \gamma v^\pi(k)] \right] \right] \\
 &= \sum_{j \in \mathcal{T}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a r_{ij}^a + \sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a r_{ij}^a + \gamma \sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a \sum_{k \in \mathcal{N} \cup \mathcal{T}} \sum_{a' \in \mathcal{A}} \pi_j^{a'} p_{jk}^{a'} r_{jk}^{a'} \\
 &\quad + \gamma^2 \sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a \sum_{k \in \mathcal{N} \cup \mathcal{T}} \sum_{a' \in \mathcal{A}} \pi_j^{a'} p_{jk}^{a'} v^\pi(k) \\
 &= \sum_{j \in \mathcal{T}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a r_{ij}^a + \sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a r_{ij}^a \\
 &\quad + \gamma \sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a \sum_{k \in \mathcal{N}} \sum_{a' \in \mathcal{A}} \pi_j^{a'} p_{jk}^{a'} r_{jk}^{a'} + \gamma \sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a \sum_{k \in \mathcal{T}} \sum_{a' \in \mathcal{A}} \pi_j^{a'} p_{jk}^{a'} r_{jk}^{a'} \\
 &\quad + \gamma^2 \sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a \sum_{k \in \mathcal{N}} \sum_{a' \in \mathcal{A}} \pi_j^{a'} p_{jk}^{a'} v^\pi(k) && \text{splitting } \mathcal{N} \text{ and } \mathcal{T}
 \end{aligned}$$

Similar to earlier, we have vectors, \mathbf{h} and \mathbf{m} with, $[\mathbf{h}]_i = \sum_{j \in \mathcal{T}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a r_{ij}^a$, $[\mathbf{m}]_i = \sum_{j \in \mathcal{N}} \sum_{a \in \mathcal{A}} \pi_i^a p_{ij}^a r_{ij}^a$, and Q is the true transition matrix of the Markov reward process induced by π and P , i.e., $[Q]_{ij} = \sum_a \pi_i^a p_{ij}^a$. The terms are not overloaded since the expectation over the true policy yields the same values. Then continuing from above, we have

$$v^\pi(i) = [\mathbf{h}]_i + [\mathbf{m}]_i + \gamma Q[\mathbf{h}]_i + \gamma Q[\mathbf{m}]_i + \gamma^2 Q^2[\mathbf{h}]_i + \gamma^2 Q^2[\mathbf{m}]_i + \dots \quad \text{unrolling } v^\pi(\mathcal{N} \cup \mathcal{T}) \quad (10)$$

$$= \left[\sum_{k=0}^{\infty} (\gamma Q)^k (\mathbf{m} + \mathbf{h}) \right]_i \quad (11)$$

$$v^\pi(i) = [(I - \gamma Q)^{-1}(\mathbf{m} + \mathbf{h})]_i \quad (12)$$

The existence of the limit and inverse are assured by Theorem A.1 in Sutton (1988). The theorem is applicable here since $\lim_{k \rightarrow \infty} (\gamma Q)^k = 0$.

B.4. MDP Certainty-Equivalence Fixed-Point

Similar to the above subsection, for certainty-equivalence fixed-point, we consider a batch of data, \mathcal{D} , with the maximum-likelihood estimate (MLE) of the above quantities according to \mathcal{D} given with a hat ($\hat{\cdot}$) on top of the quantity. The observed sets of non-terminal and terminal states in the batch are given by $\hat{\mathcal{N}}$ and $\hat{\mathcal{T}}$ respectively.

Then similar to above, we can derive the certainty-equivalence estimate of the value function according to the MLE of the policy and transition dynamics from the batch for a particular state i , $\forall i \in \hat{\mathcal{N}}$ is:

$$v^{\hat{\pi}}(i) = \left[(I - \gamma \hat{Q})^{-1}(\hat{\mathbf{m}} + \hat{\mathbf{h}}) \right]_i \quad (13)$$

This fixed-point is called the certainty-equivalence estimate (CEE) (Sutton, 1988) for an MDP. We note that MLE of the policy and transition dynamics according to the batch may not be representative of the true policy and transition dynamics. In that case, MDP-CEE (Equation (13)) is inaccurate with respect to Equation (12) due to policy and transition dynamics *sampling error*.

B.5. Policy Sampling Error Corrected MDP Certainty-Equivalence Fixed-Point

We now derive a new fixed-point, the *policy sampling error corrected MDP certainty-equivalence fixed-point*. This fixed-point corrects the policy sampling error that occurs in the value function given by Equation (13), making the estimation more accurate with respect to the true value function given by Equation (12).

We introduce the PSEC weight, $\hat{\rho}_i^a = \frac{\pi_i^a}{\hat{\pi}_i^a}$, with π being the policy that we are interested in evaluating and $\hat{\pi}$ being the MLE of the policy according to batch \mathcal{D} . $\hat{\rho}$ is then applied to the above quantities to introduce a slightly modified notation. In particular, $\hat{\rho}$ applied to \hat{Q} results in $[\hat{U}]_{ij} = \sum_{a \in \hat{\mathcal{A}}_i} \hat{\rho}_i^a \hat{\pi}_i^a \hat{p}_{ij}^a$, and applied to vectors $\hat{\mathbf{h}}$ and $\hat{\mathbf{m}}$ results in $[\hat{\mathbf{l}}]_i = \sum_{j \in \mathcal{T}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{\rho}_i^a \hat{\pi}_i^a \hat{p}_{ij}^a \bar{r}_{ij}^a$ and $[\hat{\mathbf{o}}]_i = \sum_{j \in \mathcal{N}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{\rho}_i^a \hat{\pi}_i^a \hat{p}_{ij}^a \bar{r}_{ij}^a$ respectively. After simplification, we have $[\hat{U}]_{ij} = \sum_{a \in \hat{\mathcal{A}}_i} \pi_i^a \hat{p}_{ij}^a$, $[\hat{\mathbf{l}}]_i = \sum_{j \in \mathcal{T}} \sum_{a \in \hat{\mathcal{A}}_i} \pi_i^a \hat{p}_{ij}^a \bar{r}_{ij}^a$, and $[\hat{\mathbf{o}}]_i = \sum_{j \in \mathcal{N}} \sum_{a \in \hat{\mathcal{A}}_i} \pi_i^a \hat{p}_{ij}^a \bar{r}_{ij}^a$. Using these policy sampling error corrected quantities, we can derive the fixed-point for true policy, π , in a similar manner as earlier:

$$v^\pi(i) = \left[(I - \gamma \hat{U})^{-1} (\hat{\mathbf{o}} + \hat{\mathbf{l}}) \right]_i \quad (14)$$

In computing this new fixed-point, we have corrected for the policy sampling error, resulting in a more accurate estimation of Equation (12) than Equation (13). Now, the value function is computed for the true policy that we are interested in evaluating, π .

C. Convergence of Batch Linear TD(0) to the MRP CE Fixed-Point

Theorem 1 (Batch Linear TD(0) Convergence). *For any batch whose observation vectors $\{\mathbf{x}(s) | s \in \hat{\mathcal{S}}\}$ are linearly independent, there exists an $\epsilon > 0$ such that, for all positive $\alpha < \epsilon$ and for any initial weight vector, the predictions for linear TD(0) converge under repeated presentations of the batch with weight updates after each complete presentation to the fixed-point (3).*

Proof. Batch linear TD(0) makes an update to the weight vector, \mathbf{w}_n (of dimension, length of the feature vector), after each presentation of the batch:

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \sum_{\tau \in \mathcal{D}} \sum_{t=1}^{L_\tau} \alpha \left[(\bar{r}_t + \gamma \mathbf{w}_n^T \mathbf{x}_{t+1}) - \mathbf{w}_n^T \mathbf{x}_t \right] \mathbf{x}_t$$

where \mathcal{D} is the batch of episodes, L_τ is the length of each episode τ , and α is the learning rate.

We can re-write the whole presentation of the batch of data in terms of the number of times there was a transition from state i to state j in the batch i.e. $\hat{c}_{ij} = \hat{d}_i \hat{p}_{ij}$, where \hat{d}_i is the number of times state $i \in \hat{\mathcal{N}}$ appears in the batch.

$$\begin{aligned}
 \mathbf{w}_{n+1} &= \mathbf{w}_n + \sum_{\tau \in \mathcal{D}} \sum_{t=1}^{L_\tau} \alpha [(\bar{r}_t + \gamma \mathbf{w}_n^T \mathbf{x}_{t+1} - \mathbf{w}_n^T \mathbf{x}_t)] \mathbf{x}_t \\
 &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \hat{c}_{ij} [(\bar{r}_{ij} + \gamma \mathbf{w}_n^T \mathbf{x}_j - \mathbf{w}_n^T \mathbf{x}_i)] \mathbf{x}_i \\
 &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \hat{d}_i \hat{p}_{ij} [(\bar{r}_{ij} + \gamma \mathbf{w}_n^T \mathbf{x}_j - \mathbf{w}_n^T \mathbf{x}_i)] \mathbf{x}_i \\
 &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \hat{d}_i \hat{p}_{ij} (\bar{r}_{ij} + \gamma \mathbf{w}_n^T \mathbf{x}_j) \mathbf{x}_i - \alpha \sum_{i \in \hat{\mathcal{N}}} \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \hat{d}_i \hat{p}_{ij} (\mathbf{w}_n^T \mathbf{x}_i) \mathbf{x}_i \\
 &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \hat{p}_{ij} (\bar{r}_{ij} + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i (\mathbf{w}_n^T \mathbf{x}_i) \mathbf{x}_i \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \hat{p}_{ij} \\
 &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \hat{p}_{ij} (\bar{r}_{ij} + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \mathbf{w}_n^T \mathbf{x}_i \right] && \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \hat{p}_{ij} = 1 \\
 &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\left(\sum_{j \in \hat{\mathcal{N}}} \hat{p}_{ij} (\bar{r}_{ij} + \gamma \mathbf{w}_n^T \mathbf{x}_j) \right) + \left(\sum_{j \in \hat{\mathcal{T}}} \hat{p}_{ij} \bar{r}_{ij} \right) - \mathbf{w}_n^T \mathbf{x}_i \right] && \text{If } \mathbf{x}_j \in \hat{\mathcal{T}}, \mathbf{w}_n^T \mathbf{x}_j = 0 \\
 &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\left(\sum_{j \in \hat{\mathcal{N}}} \hat{p}_{ij} \bar{r}_{ij} \right) + \left(\gamma \sum_{j \in \hat{\mathcal{N}}} \hat{p}_{ij} \mathbf{w}_n^T \mathbf{x}_j \right) + \left(\sum_{j \in \hat{\mathcal{T}}} \hat{p}_{ij} \bar{r}_{ij} \right) - \mathbf{w}_n^T \mathbf{x}_i \right] \\
 \mathbf{w}_{n+1} &= \mathbf{w}_n + \alpha \hat{X} \hat{D} \left[\hat{\mathbf{m}} + \gamma \hat{Q} \hat{X}^T \mathbf{w}_n + \hat{\mathbf{h}} - \hat{X}^T \mathbf{w}_n \right] \tag{15}
 \end{aligned}$$

where \hat{X} denotes the matrix (of dimensions, length of the feature vector by $|\hat{\mathcal{S}}|$) with columns, $\mathbf{x}_i \in \hat{\mathcal{S}}$ and \hat{D} is a diagonal matrix (of dimensions, $|\hat{\mathcal{S}}|$ by $|\hat{\mathcal{S}}|$) with $\hat{D}_{ii} = \hat{d}_i$. Given the successive updates to the weight vector \mathbf{w}_n , we now consider the actual values predicted as the following by multiplying \hat{X}^T on both sides:

$$\begin{aligned}
 \hat{X}^T \mathbf{w}_{n+1} &= \hat{X}^T \mathbf{w}_n + \alpha \hat{X}^T \hat{X} \hat{D} \left(\hat{\mathbf{m}} + \hat{\mathbf{h}} + \gamma \hat{Q} \hat{X}^T \mathbf{w}_n - \hat{X}^T \mathbf{w}_n \right) \\
 &= \hat{X}^T \mathbf{w}_n + \alpha \hat{X}^T \hat{X} \hat{D} \left(\hat{\mathbf{m}} + \hat{\mathbf{h}} \right) + \alpha \hat{X}^T \hat{X} \hat{D} \left(\gamma \hat{Q} \hat{X}^T \mathbf{w}_n - \hat{X}^T \mathbf{w}_n \right) \\
 &= \alpha \hat{X}^T \hat{X} \hat{D} \left(\hat{\mathbf{m}} + \hat{\mathbf{h}} \right) + \left(I - \alpha \hat{X}^T \hat{X} \hat{D} \left(I - \gamma \hat{Q} \right) \right) \hat{X}^T \mathbf{w}_n
 \end{aligned}$$

We then unroll the above equation by recursively applying $\hat{X}^T \mathbf{w}_n$ till $n = 0$.

$$\begin{aligned}
 \hat{X}^T \mathbf{w}_{n+1} &= \alpha \hat{X}^T \hat{X} \hat{D} \left(\hat{\mathbf{m}} + \hat{\mathbf{h}} \right) + \left(I - \alpha \hat{X}^T \hat{X} \hat{D} \left(I - \gamma \hat{Q} \right) \right) \alpha \hat{X}^T \hat{X} \hat{D} \left(\hat{\mathbf{m}} + \hat{\mathbf{h}} \right) \\
 &\quad + \left(I - \alpha \hat{X}^T \hat{X} \hat{D} \left(I - \gamma \hat{Q} \right) \right)^2 \hat{X}^T \mathbf{w}_{n-1} \\
 &\quad \vdots \\
 &= \sum_{k=0}^{n-1} \left(I - \alpha \hat{X}^T \hat{X} \hat{D} \left(I - \gamma \hat{Q} \right) \right)^k \alpha \hat{X}^T \hat{X} \hat{D} \left(\hat{\mathbf{m}} + \hat{\mathbf{h}} \right) \\
 &\quad + \left(I - \alpha \hat{X}^T \hat{X} \hat{D} \left(I - \gamma \hat{Q} \right) \right)^n \hat{X}^T \mathbf{w}_0 \tag{16}
 \end{aligned}$$

Assuming that as $n \rightarrow \infty$, $(I - \alpha \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q}))^n \rightarrow 0$, we can drop the second term and the sequence $\{\widehat{X}^T \mathbf{w}_n\}$ converges to:

$$\begin{aligned} \lim_{n \rightarrow \infty} \widehat{X}^T \mathbf{w}_n &= \left(I - (I - \alpha \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})) \right)^{-1} (\alpha \widehat{X}^T \widehat{X} \widehat{D}(\widehat{\mathbf{m}} + \widehat{\mathbf{h}})) \\ &= (I - \gamma \widehat{Q})^{-1} \widehat{D}^{-1} (\widehat{X}^T \widehat{X})^{-1} \alpha^{-1} \alpha \widehat{X}^T \widehat{X} \widehat{D}(\widehat{\mathbf{m}} + \widehat{\mathbf{h}}) \\ &= (I - \gamma \widehat{Q})^{-1} (\widehat{\mathbf{m}} + \widehat{\mathbf{h}}) \\ \lim_{n \rightarrow \infty} \mathbb{E} [\mathbf{x}_i^T \mathbf{w}_n] &= \left[(I - \gamma \widehat{Q})^{-1} (\widehat{\mathbf{m}} + \widehat{\mathbf{h}}) \right]_i, \forall i \in \widehat{\mathcal{N}} \end{aligned}$$

What is left to show now is $n \rightarrow \infty$, $(I - \alpha \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q}))^n \rightarrow 0$. Following Sutton (1988), we first show that $\widehat{D}(I - \gamma \widehat{Q})$ is positive definite, and then that $\widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})$ has a full set of eigenvalues all of whose real parts are positive. This enables us to show that α can be chosen so that eigenvalues of $(I - \alpha \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q}))$ are less than 1 in modulus, which assures us that its powers converge to 0.

To show that $\widehat{D}(I - \gamma \widehat{Q})$ is positive definite, we refer to the Gershgorin Circle theorem (Gerschgorin, 1931), which states that if a matrix, A , is real, symmetric, and strictly diagonally dominant with positive diagonal entries, then A is positive definite. However, we cannot apply this theorem as is to $\widehat{D}(I - \gamma \widehat{Q})$ since the matrix is not necessarily symmetric. To use the theorem, we first apply another theorem (Theorem A.3 from Sutton (1988)) that states: a square matrix A is positive definite if and only if $A + A^T$ is positive definite. So it suffices to show that $\widehat{D}(I - \gamma \widehat{Q}) + (\widehat{D}(I - \gamma \widehat{Q}))^T$ is positive definite.

Consider the matrix $S = \widehat{D}(I - \gamma \widehat{Q}) + (\widehat{D}(I - \gamma \widehat{Q}))^T$. We know that S is real and symmetric. It remains to show that the diagonal entries are positive and that S is strictly diagonally dominant. First, we look at the diagonal entries, $S_{ii} = 2[\widehat{D}(I - \gamma \widehat{Q})]_{ii} = 2\hat{d}_i(1 - \gamma \hat{p}_{ii}) > 0, \forall i \in \widehat{\mathcal{N}}$, which are positive. Second, we have the non-diagonal entries for $i \neq j$ as $S_{ij} = [\widehat{D}(I - \gamma \widehat{Q})]_{ij} + [\widehat{D}(I - \gamma \widehat{Q})]_{ji} = -\gamma \hat{d}_i \hat{p}_{ij} - \gamma \hat{d}_j \hat{p}_{ji} \leq 0$, which are nonpositive. We want to show that $|S_{ii}| \geq \sum_{j \neq i} |S_{ij}|$, with strict inequality holding for at least one i ; we know that the diagonal elements $S_{ii} > 0$ and non-diagonal elements $S_{ij} \leq 0, i \neq j$. Hence, to show that S is strictly diagonally dominant, it is enough to show that $S_{ii} > -\sum_{j \neq i} S_{ij}$, which means we can simply show that the sum of each entire row is greater than 0, i.e. $\sum_j S_{ij} > 0$.

Before we show that $\sum_j S_{ij} > 0$, we note that $\tilde{d}^T = \hat{\mu}^T (I - \widehat{Q})^{-1}$ where $\hat{\mu}_i$ is the empirical state distribution of state i . Given the definitions of \hat{d} , $\hat{\mu}$, and \widehat{Q} , this fact follows from Kemeny et al. (1960) and is used by Sutton (1988). Using this fact, we show that $\sum_j S_{ij} \geq 0$:

$$\begin{aligned} \sum_j S_{ij} &= \sum_j \left([\widehat{D}(I - \gamma \widehat{Q})]_{ij} + [\widehat{D}(I - \gamma \widehat{Q})]_{ji}^T \right) \\ &= \hat{d}_i \sum_j \left([I - \gamma \widehat{Q}]_{ij} + \sum_j \hat{d}_j [I - \gamma \widehat{Q}]_{ij}^T \right) \\ &= \hat{d}_i \sum_j (1 - \gamma \hat{p}_{ij}) + \left[\tilde{d}^T (I - \widehat{Q}) \right]_i \\ &= \hat{d}_i \sum_j (1 - \gamma \hat{p}_{ij}) + \left[\hat{\mu}^T (I - \widehat{Q})^{-1} (I - \widehat{Q}) \right]_i && \tilde{d}^T = \hat{\mu}^T (I - \widehat{Q})^{-1} \\ &= \hat{d}_i (1 - \gamma \sum_j \hat{p}_{ij}) + \hat{\mu}_i \\ &\geq 0, \end{aligned}$$

where the final inequality is strict since $\hat{\mu}$ is positive for at least one element. Given the above, we have shown that S is real, symmetric, and strictly diagonally dominant; hence, S is positive definite according to the Gershgorin Circle theorem (Gerschgorin, 1931). Since, $S = \widehat{D}(I - \gamma \widehat{Q}) + (\widehat{D}(I - \gamma \widehat{Q}))^T$ is positive definite, we have $\widehat{D}(I - \gamma \widehat{Q})$ to be positive definite.

Now we need to show that $\widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})$ has a full set of eigenvalues, all of whose real parts are positive. We know that $\widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})$ has a full set of eigenvalues for the same reason shown by Sutton (1988), i.e. $\widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})$ is a

product of three non-singular matrices, which means $\widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})$ is nonsingular as well; hence, no eigenvalues are 0 i.e. its set of eigenvalues is full.

Consider λ and y to be an eigenvalue and eigenvector pair of $\widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})$. First lets consider that y may be a complex number and is of the form $y = a + bi$, and let $z = (\widehat{X}^T \widehat{X})^{-1}y, y \neq 0$. Second, we consider $\widehat{D}(I - \gamma \widehat{Q})$ from earlier i.e. where $*$ is the conjugate-transpose:

$$\begin{aligned}
 y^* \widehat{D}(I - \gamma \widehat{Q})y &= z^* \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})y && \text{substituting } y^* \\
 &= z^* \lambda y && \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q})y = \lambda y \\
 &= \lambda z^* \widehat{X}^T \widehat{X} z && \text{substituting } y \\
 &= \lambda (\widehat{X} z)^* \widehat{X} z && \\
 (a^T - b^T i)(\widehat{D}(I - \gamma \widehat{Q}))(a^T + b^T i) &= \lambda (\widehat{X} z)^* \widehat{X} z && \text{substituting } y^* \text{ and } y
 \end{aligned}$$

From the above equality, we know that the real parts (Re) of the LHS and RHS are equal as well i.e.

$$\begin{aligned}
 \text{Re} \left(y^* \widehat{D}(I - \gamma \widehat{Q})y \right) &= \text{Re} \left(\lambda (\widehat{X} z)^* \widehat{X} z \right) \\
 a^T \widehat{D}(I - \gamma \widehat{Q})a + b^T \widehat{D}(I - \gamma \widehat{Q})b &= (\widehat{X} z)^* \widehat{X} z \text{Re}(\lambda)
 \end{aligned}$$

LHS must be strictly positive since we already proved that $\widehat{D}(I - \gamma \widehat{Q})$ is positive definite and by definition, RHS, $(\widehat{X} z)^* \widehat{X} z$, is strictly positive as well. Thus, the $\text{Re}(\lambda)$ must be positive. Finally, using this result we want to show that the eigenvalues of $(I - \alpha \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q}))$ are of modulus less than 1 for a suitable α .

First, we can see that y is also an eigenvector of $(I - \alpha \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q}))$, since $(I - \alpha \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q}))y = y - \alpha \lambda y = (1 - \lambda \alpha)y$, where $\lambda' = (1 - \alpha \lambda)$ is an eigenvalue of $(I - \alpha \widehat{X}^T \widehat{X} \widehat{D}(I - \gamma \widehat{Q}))$. Second, we want to find suitable α such that the modulus of λ' is less than 1. We have the modulus of λ' :

$$\begin{aligned}
 \|\lambda'\| &= \|1 - \alpha \lambda\| \\
 &= \sqrt{(1 - \alpha a)^2 + (-\alpha b)^2} && \text{substituting } \lambda = a + bi \text{ of general complex form} \\
 &= \sqrt{1 - 2\alpha a + \alpha^2 a^2 + \alpha^2 b^2} \\
 &= \sqrt{1 - 2\alpha a + \alpha^2 (a^2 + b^2)} \\
 &< \sqrt{1 - 2\alpha a + \alpha \frac{2a}{(a^2 + b^2)} (a^2 + b^2)} && \text{using } \alpha = \frac{2a}{(a^2 + b^2)} \\
 &= \sqrt{1 - 2\alpha a + 2\alpha a} = 1
 \end{aligned}$$

From above, we can see that if α is chosen such that $0 < \alpha < \frac{2a}{a^2 + b^2}$, then λ' will have modulus less than 1. Then using the theorem that states: if a matrix A has n independent eigenvectors with eigenvalues λ_i , then $A^k \rightarrow 0$ as $k \rightarrow \infty$ if and only if all $\|\lambda_i\| < 1$, which implies that $\lim_{n \rightarrow \infty} \left(I - \alpha \widehat{X}^T \widehat{D}(I - \widehat{Q}) \widehat{X}^T \right)^n = 0$, taking the trailing element in Equation (16) to 0 for a suitable α . We thus prove convergence to the fixed point in Equation (3) if a batch linear TD(0) update is used with an appropriate step size α . \square

D. Convergence of Batch Linear TD(0) to the MDP CE Fixed-Point

Theorem 2 (Batch Linear TD(0) Convergence). *For any batch whose observation vectors $\{\mathbf{x}(s) | s \in \widehat{\mathcal{S}}\}$ are linearly independent, there exists an $\epsilon > 0$ such that, for all positive $\alpha < \epsilon$ and for any initial weight vector, the predictions for linear TD(0) converge under repeated presentations of the batch with weight updates after each complete presentation to the fixed-point (4).*

Proof. Batch linear TD(0) makes an update to weight vector, \mathbf{w}_n (of dimension, length of the feature vector), after each presentation of the batch:

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \sum_{\tau \in \mathcal{D}} \sum_{t=1}^{L_\tau} \alpha [(\bar{r}_t + \gamma \mathbf{w}_n^T \mathbf{x}_{t+1}) - \mathbf{w}_n^T \mathbf{x}_t] \mathbf{x}_t$$

where \mathcal{D} is the batch of episodes, L_τ is the length of each episode τ , and α is the learning rate.

We can re-write the whole presentation of the batch of data in terms of the number of times there was a transition from state i to state j when taking action a in the batch i.e. $\hat{c}_{ij}^a = \hat{d}_i \hat{\pi}_i^a \hat{p}_{ij}^a$, where \hat{d}_i is the number of times state $i \in \hat{\mathcal{N}}$ appears in the batch.

$$\begin{aligned} \mathbf{w}_{n+1} &= \mathbf{w}_n + \sum_{\tau \in \mathcal{D}} \sum_{t=1}^{L_\tau} \alpha [(\bar{r}_t + \gamma \mathbf{w}_n^T \mathbf{x}_{t+1} - \mathbf{w}_n^T \mathbf{x}_t)] \mathbf{x}_t \\ &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{c}_{ij}^a [(\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j - \mathbf{w}_n^T \mathbf{x}_i)] \mathbf{x}_i \\ &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{d}_i \hat{\pi}_i^a \hat{p}_{ij}^a [(\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j - \mathbf{w}_n^T \mathbf{x}_i)] \mathbf{x}_i \\ &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{d}_i \hat{p}_{ij}^a \hat{\pi}_i^a (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) \mathbf{x}_i - \alpha \sum_{i \in \hat{\mathcal{N}}} \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{d}_i \hat{\pi}_i^a \hat{p}_{ij}^a (\mathbf{w}_n^T \mathbf{x}_i) \mathbf{x}_i \\ &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{p}_{ij}^a \hat{\pi}_i^a (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i (\mathbf{w}_n^T \mathbf{x}_i) \mathbf{x}_i \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{\pi}_i^a \hat{p}_{ij}^a \\ &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{p}_{ij}^a \hat{\pi}_i^a (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \mathbf{w}_n^T \mathbf{x}_i \right] \sum_{j \in \hat{\mathcal{N}} \cup \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{\pi}_i^a \hat{p}_{ij}^a = 1 \\ &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\left(\sum_{j \in \hat{\mathcal{N}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{p}_{ij}^a \hat{\pi}_i^a (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) \right) + \left(\sum_{j \in \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{p}_{ij}^a \hat{\pi}_i^a \bar{r}_{ij}^a \right) - \mathbf{w}_n^T \mathbf{x}_i \right] \text{ If } \mathbf{x}_j \in \hat{\mathcal{T}}, \mathbf{w}_n^T \mathbf{x}_j = 0 \\ &= \mathbf{w}_n + \alpha \sum_{i \in \hat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\left(\sum_{j \in \hat{\mathcal{N}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{p}_{ij}^a \hat{\pi}_i^a \bar{r}_{ij}^a \right) + \left(\gamma \sum_{j \in \hat{\mathcal{N}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{p}_{ij}^a \hat{\pi}_i^a \mathbf{w}_n^T \mathbf{x}_j \right) + \left(\sum_{j \in \hat{\mathcal{T}}} \sum_{a \in \hat{\mathcal{A}}_i} \hat{p}_{ij}^a \hat{\pi}_i^a \bar{r}_{ij}^a \right) - \mathbf{w}_n^T \mathbf{x}_i \right] \\ &\mathbf{w}_{n+1} = \mathbf{w}_n + \alpha \hat{X} \hat{D} \left[\hat{\mathbf{m}} + \gamma \hat{Q} \hat{X}^T \mathbf{w}_n + \hat{\mathbf{h}} - \hat{X}^T \mathbf{w}_n \right] \end{aligned} \quad (17)$$

where \hat{X} denotes the matrix (of dimensions, length of the feature vector by $|\hat{\mathcal{S}}|$) with columns, $\mathbf{x}_i \in \hat{\mathcal{S}}$ and \hat{D} is a diagonal matrix (of dimensions, $|\hat{\mathcal{S}}|$ by $|\hat{\mathcal{S}}|$) with $\hat{D}_{ii} = \hat{d}_i$.

Notice that Equation (17) is the same as Equation (15) since the considered MRP and MDP settings are equivalent. Due to this similarity, we omit the proof from here below as it is identical to the Theorem 1 proof. \square

E. Convergence of Batch Linear PSEC-TD(0) to the PSEC-MDP-CE Fixed-Point

E.1. PSEC Correction Applied to the New Estimate

We now show that batch linear PSEC-TD(0) converges to the policy corrected MDP-CE established in Equation (14).

Theorem 3 (Batch Linear PSEC-TD(0) Convergence). *For any batch whose observation vectors $\{\mathbf{x}(s)|s \in \widehat{\mathcal{S}}\}$ are linearly independent, there exists an $\epsilon > 0$ such that, for all positive $\alpha < \epsilon$ and for any initial weight vector, the predictions for linear PSEC-TD(0) converge under repeated presentations of the batch with weight updates after each complete presentation to the fixed-point (5).*

Proof. The proof for PSEC-TD(0) follows in large part the structure of the proof for TD(0). Below we highlight the salient points in the proof.

Batch linear PSEC-TD(0) makes an update to the weight vector, \mathbf{w}_n (of dimension, length of the feature vector), after each presentation of the batch:

$$\mathbf{w}_{n+1} = \mathbf{w}_n + \sum_{\tau \in \mathcal{D}} \sum_{t=1}^{L_\tau} \alpha [\hat{\rho}_t(\bar{r}_t + \gamma \mathbf{w}_n^T \mathbf{x}_{t+1}) - \mathbf{w}_n^T \mathbf{x}_t] \mathbf{x}_t$$

where \mathcal{D} is the batch of episodes, L_τ is the length of each episode τ , $\hat{\rho}_t$ is the PSEC correction weight at time t for a given episode τ , and α is the learning rate.

We can re-write the whole presentation of the batch of data in terms of the number of times there was a transition from state i to state j when taking action a in the batch i.e. $\hat{c}_{ij}^a = \hat{d}_i \hat{\pi}_i^a \hat{p}_{ij}^a$, where \hat{d}_i is the number of times state $i \in \widehat{\mathcal{N}}$ appears in the batch.

$$\begin{aligned} \mathbf{w}_{n+1} &= \mathbf{w}_n + \sum_{\tau \in \mathcal{D}} \sum_{t=1}^{L_\tau} \alpha [\hat{\rho}_t(\bar{r}_t + \gamma \mathbf{w}_n^T \mathbf{x}_{t+1}) - \mathbf{w}_n^T \mathbf{x}_t] \mathbf{x}_t \\ &= \mathbf{w}_n + \alpha \sum_{i \in \widehat{\mathcal{N}}} \sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{c}_{ij}^a [\hat{\rho}_i^a(\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \mathbf{w}_n^T \mathbf{x}_i] \mathbf{x}_i \\ &= \mathbf{w}_n + \alpha \sum_{i \in \widehat{\mathcal{N}}} \sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{d}_i \hat{\pi}_i^a \hat{p}_{ij}^a [\hat{\rho}_i^a(\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \mathbf{w}_n^T \mathbf{x}_i] \mathbf{x}_i \\ &= \mathbf{w}_n + \alpha \sum_{i \in \widehat{\mathcal{N}}} \sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{d}_i \hat{\pi}_i^a \hat{p}_{ij}^a \left[\frac{\pi_i^a}{\hat{\pi}_i^a} (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \mathbf{w}_n^T \mathbf{x}_i \right] \mathbf{x}_i \\ &= \mathbf{w}_n + \alpha \sum_{i \in \widehat{\mathcal{N}}} \sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{d}_i \hat{p}_{ij}^a \pi_i^a (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) \mathbf{x}_i - \alpha \sum_{i \in \widehat{\mathcal{N}}} \sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{d}_i \hat{\pi}_i^a \hat{p}_{ij}^a (\mathbf{w}_n^T \mathbf{x}_i) \mathbf{x}_i \\ &= \mathbf{w}_n + \alpha \sum_{i \in \widehat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{p}_{ij}^a \pi_i^a (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \alpha \sum_{i \in \widehat{\mathcal{N}}} \hat{d}_i (\mathbf{w}_n^T \mathbf{x}_i) \mathbf{x}_i \sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{\pi}_i^a \hat{p}_{ij}^a \\ &= \mathbf{w}_n + \alpha \sum_{i \in \widehat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{p}_{ij}^a \pi_i^a (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) - \mathbf{w}_n^T \mathbf{x}_i \right] \sum_{j \in \widehat{\mathcal{N}} \cup \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{\pi}_i^a \hat{p}_{ij}^a = 1 \\ &= \mathbf{w}_n + \alpha \sum_{i \in \widehat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\left(\sum_{j \in \widehat{\mathcal{N}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{p}_{ij}^a \pi_i^a (\bar{r}_{ij}^a + \gamma \mathbf{w}_n^T \mathbf{x}_j) \right) + \left(\sum_{j \in \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{p}_{ij}^a \pi_i^a \bar{r}_{ij}^a \right) - \mathbf{w}_n^T \mathbf{x}_i \right] \text{ If } \mathbf{x}_j \in \widehat{\mathcal{T}}, \mathbf{w}_n^T \mathbf{x}_j = 0 \\ &= \mathbf{w}_n + \alpha \sum_{i \in \widehat{\mathcal{N}}} \hat{d}_i \mathbf{x}_i \left[\left(\sum_{j \in \widehat{\mathcal{N}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{p}_{ij}^a \pi_i^a \bar{r}_{ij}^a \right) + \left(\gamma \sum_{j \in \widehat{\mathcal{N}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{p}_{ij}^a \pi_i^a \mathbf{w}_n^T \mathbf{x}_j \right) + \left(\sum_{j \in \widehat{\mathcal{T}}} \sum_{a \in \widehat{\mathcal{A}}_i} \hat{p}_{ij}^a \pi_i^a \bar{r}_{ij}^a \right) - \mathbf{w}_n^T \mathbf{x}_i \right] \\ &= \mathbf{w}_n + \alpha \widehat{X} \widehat{D} \left[\widehat{\mathbf{o}} + \gamma \widehat{U} \widehat{X}^T \mathbf{w}_n + \widehat{\mathbf{1}} - \widehat{X}^T \mathbf{w}_n \right] \end{aligned}$$

where \widehat{X} denotes the matrix (of dimensions, length of the feature vector by $|\widehat{\mathcal{S}}|$) with columns, $\mathbf{x}_i \in \widehat{\mathcal{S}}$ and \widehat{D} is a diagonal matrix (of dimensions, $|\widehat{\mathcal{S}}|$ by $|\widehat{\mathcal{S}}|$) with $\widehat{D}_{ii} = \hat{d}_i$.

Assuming that as $n \rightarrow \infty$, $(I - \alpha \hat{X}^T \hat{X} \hat{D} (I - \gamma \hat{U}))^n \rightarrow 0$, we can drop the second term and the sequence $\{\hat{X}^T \mathbf{w}_n\}$ converges to:

$$\begin{aligned} \lim_{n \rightarrow \infty} \hat{X}^T \mathbf{w}_n &= \left(I - (I - \alpha \hat{X}^T \hat{X} \hat{D} (I - \gamma \hat{U})) \right)^{-1} (\alpha \hat{X}^T \hat{X} \hat{D} (\hat{\mathbf{o}} + \hat{\mathbf{1}})) \\ &= (I - \gamma \hat{U})^{-1} \hat{D}^{-1} (\hat{X}^T \hat{X})^{-1} \alpha^{-1} \alpha \hat{X}^T \hat{X} \hat{D} (\hat{\mathbf{o}} + \hat{\mathbf{1}}) \\ &= (I - \gamma \hat{U})^{-1} (\hat{\mathbf{o}} + \hat{\mathbf{1}}) \\ \lim_{n \rightarrow \infty} \mathbb{E} [\mathbf{x}_i^T \mathbf{w}_n] &= \left[(I - \gamma \hat{U})^{-1} (\hat{\mathbf{o}} + \hat{\mathbf{1}}) \right]_i, \forall i \in \hat{\mathcal{N}} \end{aligned}$$

What is left to show now is that as $n \rightarrow \infty$, $(I - \alpha \hat{X}^T \hat{X} \hat{D} (I - \gamma \hat{U}))^n \rightarrow 0$, which we can show by following the steps shown for Equation (16). Thus we prove convergence to the fixed-point (5). \square

E.2. Convergence to the MDP True Fixed-Point with Infinite Data

With batch linear PSEC-TD(0) we have corrected for the policy sampling error in batch linear TD(0). The remaining inaccuracy of the policy sampling corrected certainty-equivalence fixed-point is due to the transition dynamics sampling error. In a model-free setting, however, we cannot correct for this error in the same way we corrected the policy sampling error.

We argue that as the batch size approaches infinite, the maximum-likelihood estimate of the transition dynamics will approach the true transition dynamics i.e. $\hat{p} \rightarrow p$. It then follows that in expectation, the true value function will be reached. Thus, the batch linear PSEC-TD(0) with an infinite batch size will correctly converge to the true value function fixed-point given by Equation (12).

F. Additional Empirical Results

F.1. Tabular Setting: Discrete States and Actions

F.1.1. OFF-POLICY RESULTS

For off-policy TD(0), we always use the variant that applies the importance weight to the TD-error.

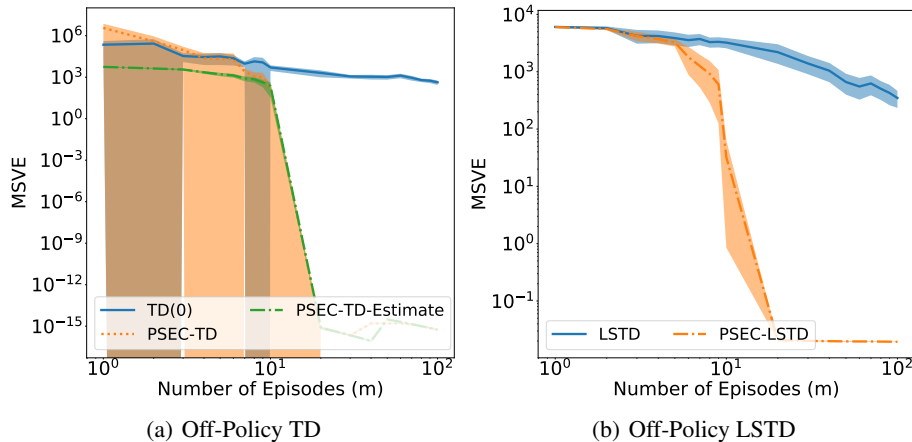


Figure 5. Deterministic Gridworld experiments. Both axes are log-scaled, leading to the asymmetric confidence intervals. Errors are computed over 50 trials with 95% confidence intervals. Figure 5(a) and Figure 5(b) compare the final errors achieved by variants of PSEC-TD(0) and TD(0), and PSEC-LSTD(0) and LSTD(0) respectively for varying batch size in the off-policy case.

For off-policy TD(0) we only consider the PSEC-TD variant as we found multiplying the new estimate by the weight was divergent. All methods show higher variance for the off-policy setting, however, PSEC-TD(0) variants still provide more

accurate value function estimates.

F.1.2. EFFECT OF ENVIRONMENT STOCHASTICITY

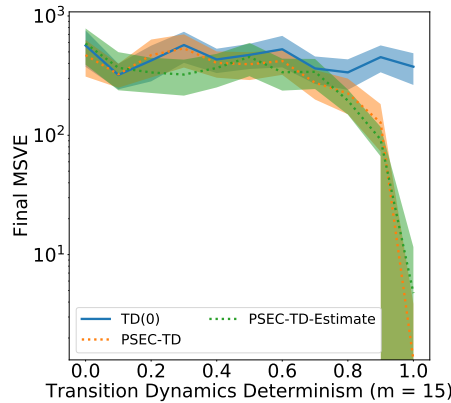


Figure 6. Additional Gridworld experiments. Errors are computed over 50 trials with 95% confidence intervals. Figure 6 is a y -axis log scaled graph that shows the final error (averaged over 100 trials) achieved by the two variants of PSEC-TD(0) and TD(0) for a given batch size (15 episodes) with varying levels of determinism of the transition dynamics.

According to Theorem 2, TD suffers from policy *and* transition dynamics sampling error. We study this observation through Figure 6, which illustrates how the performance of PSEC changes with different levels of transition dynamics determinism for a fixed batch size. In Gridworld, the determinism is varied according to a parameter, p , where the environment becomes purely deterministic or stochastic as $p \rightarrow 1$ or $p \rightarrow 0$ respectively. From Theorem 3 we expect PSEC to fully correct for the policy sampling error but not transition dynamics sampling error. Figure 6 confirms that PSEC achieves a lower final MSVE than TD as $p \rightarrow 1$. As $p \rightarrow 0$, the transition dynamics become the dominant source of sampling error and PSEC-TD(0) and TD(0) perform similarly.

F.2. Function Approximation: Continuous States and Discrete Actions (CartPole)

In each experiment below, unless stated, the following components were fixed: a batch size of 10 episodes, the value function was represented with a neural network of single hidden layer of 512 neurons using tanh activation, the gradients were normalized to unit norm before the gradient descent step was performed, we used a learning rate of 1.0 and decayed the learning rate by 10% every 50 presentations of the batch to the algorithm. The true MSVE was computed by 200 Monte Carlo rollouts for 150 sampled states.

F.2.1. DATA EFFICIENCY

From Figure 3(a). Both algorithms used a learning rate of 1.0 and decayed the learning rate by 5% every 10 presentations of the batch. The PSEC model architecture was a neural network with 3 hidden layers with 16 neurons each. Batch sizes of 10, 50, and 500 episodes used a learning rate of 0.0125, batch size of 100 used 0.003125, and batch size of 1000 used 0.001563.

F.2.2. EFFECT OF VALUE FUNCTION MODEL ARCHITECTURE

From Figure 4(a). PSEC used a model architecture of 3 hidden layers with 16 neurons each and tanh activation, and learning rate of 0.025.

F.2.3. EFFECT OF PSEC LEARNING RATE

Figure 7 compares the data efficiency of PSEC-TD vs TD for varying learning rates of the PSEC policy, while holding the value function, PSEC model, and behavior policy architectures fixed, on CartPole. Since TD does not use PSEC, its error for a given batch size is independent of the PSEC learning rate. From above, PSEC appears to be relatively stable in its improvement over TD regardless of the learning rate used. The PSEC policy used in this experiment was a neural network

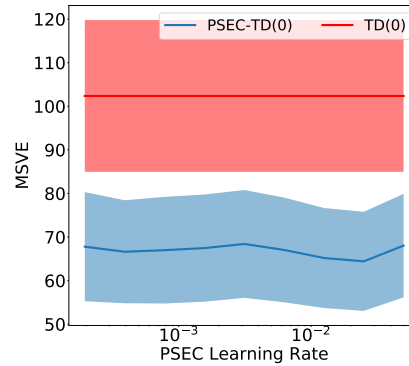


Figure 7. Compares data efficiency of PSEC with varying learning rates against TD respectively on CartPole. Results use a batch size of 10 episodes, and results are averaged over 300 trials with 95% confidence error bars.

with: 3 hidden layers with 16 neurons each and tanh activation.

F.2.4. EFFECT OF PSEC MODEL ARCHITECTURE

From Figure 4(b). PSEC used a learning rate of 0.025. The chosen PSEC neural network models are with respect to the behavior policy described earlier, a 2 hidden layered with 16 neurons architecture.

F.2.5. VARYING PSEC TRAINING STYLE

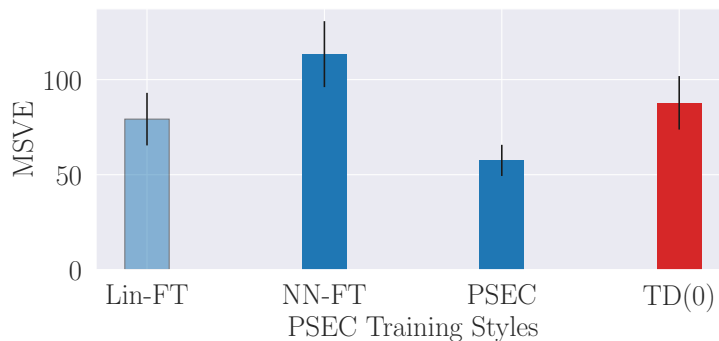


Figure 8. Comparing data efficiency of varying training styles of PSEC against TD for a fixed batch size of 10 episodes. Results shown are averaged over 300 trials and shaded region is 95% confidence. Darker shades represent statistically significant result.

Figure 8 illustrates the data efficiency of three variants of PSEC, while holding the value-function PSEC model, and behavior policy architectures fixed. All three variants use the same PSEC model architecture as that of the behavior policy, and each used a learning rate of 0.025 with tanh activation. The three variants are as follows: 1) Lin-FT is when PSEC initializes the PSEC model to the weights of the behavior policy and trains on the batch of data by finetuning only the last linear layer, 2) NN-FT is when PSEC initializes the PSEC model to the weights of the behavior policy but finetunes the all the weights of the network, and 3) PSEC uses the same training style in the previous experiments, where the model is initialized randomly and all the weights are tuned. We found that Lin-FT performed similarly to TD with a statistically insignificant improvement over TD; we believe this may be so since Lin-FT is initialized to the behavior policy and since there are only few weights to change in the linear layer, the newly learned Lin-FT is still similar to the behavior policy, which would produce PSEC corrections close to 1 (equivalent to TD). Interestingly, tuning all the weights of the neural network did better when the model was initialized randomly versus when it was initialized to the behavior policy.

F.2.6. EFFECT OF UNDERFITTING AND OVERFITTING DURING PSEC POLICY TRAINING

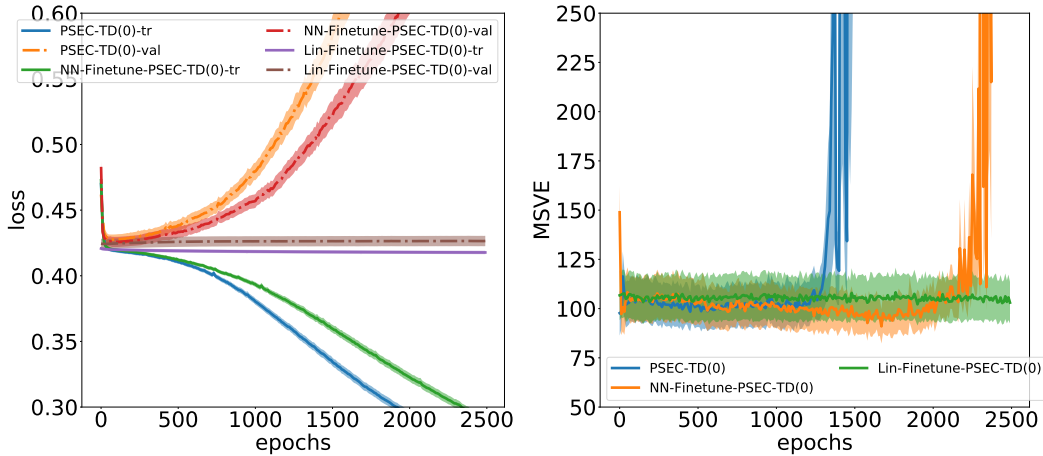


Figure 9. Comparing MSVE achieved and cross entropy loss (training (tr) and validation (val)) by variants of PSEC after each epoch of training for a fixed batch size of 10 episodes. Results shown are averaged over 50 trials and shaded region is 95% confidence.

This experiment attempts to give an understanding of how the MSVE achieved by each PSEC variant is dependent on the number of epochs the PSEC model was trained for, while holding the value-function, PSEC and behavior policy architectures fixed. We conduct the experiment as follows: the PSEC algorithm performs 10 gradient descent steps (epochs) on the full batch of data, after which the resulting training and validation mean cross-entropy losses are plotted along with the MSVE achieved by that trained PSEC policy. For example, after 10 epochs, the training and validation loss of the PSEC model was nearly 0.5 and the model achieved an MSVE of nearly 150.

Since computing the MSVE can be computationally expensive, as it requires processing the batch until the value function converges, we change the learning rate decay schedule to starting with a learning rate of 1.0 but decaying learning rate by 50% every 50 presentations of the full batch to the algorithm (this change is also the reason why these results may be different from the ones shown earlier). All PSEC variants used a learning rate of 0.025 and PSEC model architecture of 2 hidden layers with 16 neurons each and tanh activation.

Figure 9 suggests that performance of PSEC, regardless of the variant, depends on the number of epochs it was trained for. Naturally, we do want to fit sufficiently well to the data, and the graph suggests that some overfitting is tolerable. However, if overfitting becomes extreme, PSEC's performance suffers, resulting in MSVE nearly 1000 times larger than the minimum error achieved (not shown for clarity). From the graph, we can see that the PSEC variant, which is initialized randomly, starts to extremely overfit before the NN-finetune variant does, causing its MSVE to degrade before that of NN-finetune variant. We also see that the Lin-finetune variant is not able to overfit since the last linear layer may not be expressive enough to overfit, causing it to have a relatively stable MSVE across all epochs.

F.2.7. EFFECT OF BEHAVIOR POLICY DISTRIBUTION

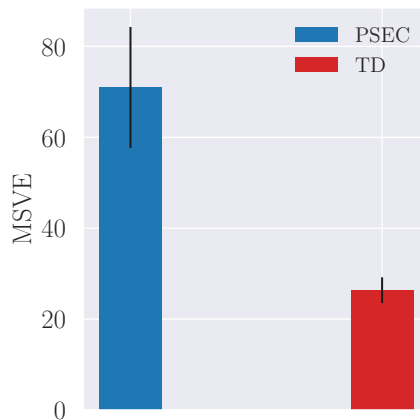


Figure 10. Comparing data efficiency of PSEC against TD for a fixed batch size of 100 episodes when the behavior policy models a discontinuous function. Results shown are averaged over 30 trials and shaded region is 95% confidence.

So far, PSEC has used a function approximator of the similar function class as that of the behavior policy i.e. both the PSEC models and the behavior policy were neural networks of similar architectures. In this experiment, we evaluate the performance of PSEC with a neural network policy when the behavior policy that models a discontinuous function generates a larger batch size of 100 episodes, while the value function and PSEC model architectures are fixed. In particular, we use a behavior policy in CartPole that does the following: if the sign of the pole angle is negative, move left with probability 0.75 and right with probability 0.25, and if the sign of the pole angle is positive, move right with probability 0.75 and left with probability 0.25. The PSEC policy is a neural network with 3 hidden layers with 16 neurons each with tanh activation.

Figure 10 shows that PSEC performs much worse than TD when the behavior policy is the discontinuous type function described above. We reason that the neural network finds it difficult to compute the MLE of the data since this discontinuous distribution is “hard” to model; therefore, producing incorrect PSEC weights, which degrade its performance. While we can largely ignore the distribution of the behavior policy, this experiment shows that PSEC may suffer in situations like the one described.

F.3. Function Approximation: Continuous States and Actions (InvertedPendulum)

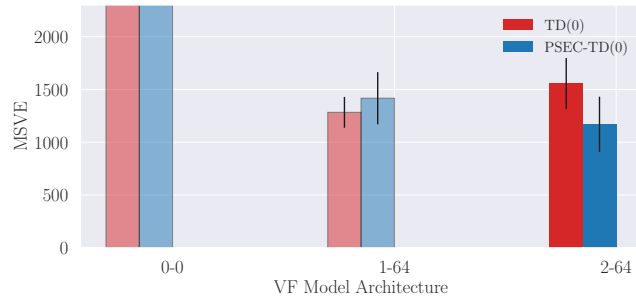
In each experiment below, unless stated, the following components were fixed: a batch size of 20 episodes, the value function was represented with a neural network of 2 hidden layer with 64 neurons each using tanh activation, the gradients were normalized to unit norm before the gradient descent step was performed, we used a learning rate of 1.0 and decayed the learning rate by 5% every 10 presentations of the batch to the algorithm. The true MSVE was computed by 100 Monte Carlo rollouts for 100 sampled states.

F.3.1. DATA EFFICIENCY

From Figure 3(b). The PSEC model architecture was a neural network with 2 hidden layers with 64 neurons each. Batch sizes 10, 50, 100, 500, 1000 used a learning rate of 0.000781. Batch size of 10 used an L2 weight penalization of 0.02. All used a value function model architecture of 3 hidden layers with 64 neurons each.

F.3.2. EFFECT OF VALUE FUNCTION MODEL ARCHITECTURE

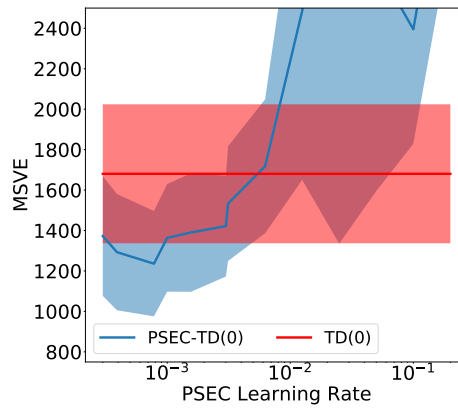
From Figure 11(a). The PSEC policy is 2 hidden layers with 64 neurons each and used a learning rate of 0.000781.



(a) InvertedPendulum

Figure 11. Comparing data efficiency of PSEC with varying VF model architectures against TD. Figure 11(a) use batch size of 20 episodes, and results shown are averaged over 350 trials respectively with error bars of 95% confidence. Darker shades represent statistically significant result. The label on the x axis shown is (# hidden layers - # neurons).

F.3.3. EFFECT OF PSEC LEARNING RATE



(a)

Figure 12. Comparing data efficiency of PSEC with varying learning rates against TD for a fixed batch size of 20 episodes. Results shown are averaged over 200 trials and error bar is 95% confidence.

Figure 12 compares the data efficiency of PSEC-TD to TD with varying learning rates for PSEC, while holding the PSEC and value function architecture fixed. Unlike earlier, the PSEC learning rate heavily influences the learned value function in the continuous state and action setting. In general, PSEC performance heavily degraded when the learning rate increased (y -axis limited for clarity). Among the tested learning rates, 0.000781 was the optimal, giving a statistically significant result.

F.3.4. EFFECT OF PSEC MODEL ARCHITECTURE

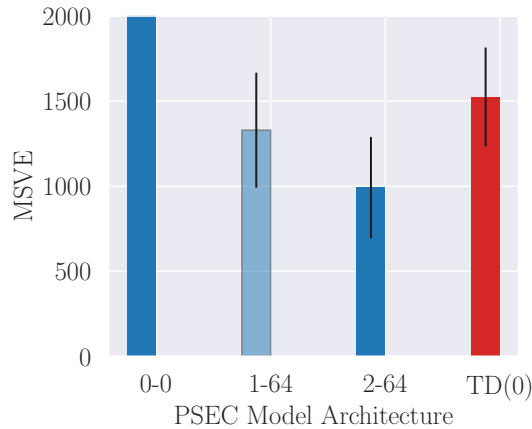


Figure 13. Comparing data efficiency of PSEC with varying model architectures against TD for a fixed batch size of 20 episodes. Results shown are averaged over 200 trials and error bar is 95% confidence. Darker shades represent statistically significant result. The label on the x axis shown is (# hidden layers - # neurons). The value function represented by 0-0 is a linear mapping with no activation function.

Figure 13 compares the data efficiency of PSEC-TD with TD with varying PSEC model architectures, while holding the value-function and behavior policy architectures fixed. All the shown PSEC architectures used a learning rate of learning rate 0.000781. Similar to our earlier findings, a more expressive network was able to better model the batch of data and produce a statistically significant improvement over TD. Less expressive PSEC models performed worse than TD, and any improvement was statistically insignificant. Note that the linear architecture used produced an MSVE of ~ 5800 (not shown for clarity) and its poor data efficiency with respect to TD(0) was statistically significant.

F.3.5. EFFECT OF UNDERFITTING AND OVERFITTING DURING PSEC POLICY TRAINING

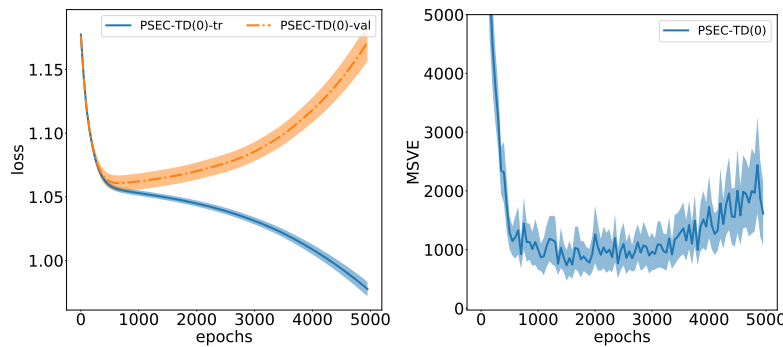


Figure 14. Comparing MSVE achieved, and training (tr) and validation (val) loss by PSEC after each epoch of training for a fixed batch size of 20 episodes on InvertedPendulum. Results shown are averaged over 100 trials and shaded region is 95% confidence.

This experiment attempts to give an understanding of how the MSVE achieved by each PSEC variant is dependent on the number of epochs the PSEC model was trained for, while holding the value-function, PSEC and behavior policy architectures fixed on InvertedPendulum. The experiment is conducted in a similar manner as before except we perform 50 gradient descent steps (epochs) before plotting the training and validation loss, and MSVE achieved by PSEC after the gradient steps. The loss shown here is a regression loss detailed in Appendix G.3. Figure 14 suggests that performance of PSEC depends on the number of epochs it was trained for. Naturally, we do want to fit sufficiently well to the data, and the graph suggests that some overfitting is tolerable. However, if overfitting becomes extreme, PSEC’s performance suffers. If some overfitting is

desirable, then early stopping is not the preferred principled approach to terminate PSEC model training. When computing MSVE, we used the learning rate schedule specified at the beginning of this section. PSEC used a learning rate of 0.000781 and model architecture of 2 hidden layers with 64 neurons each and tanh activation.

G. Extended Empirical Description

In this appendix we provide additional details for our empirical evaluation.

G.1. Gridworld

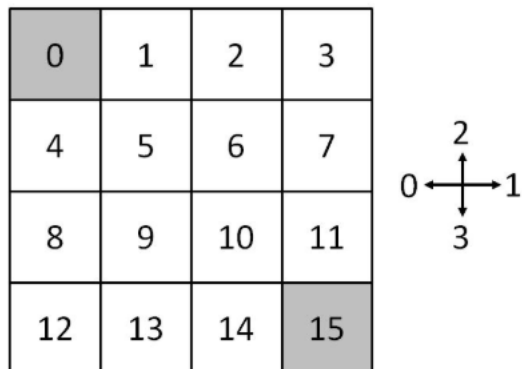


Figure 15. The Gridworld environment. Start at top left, bottom right is terminal state, discrete action space consists of the cardinal directions, and discrete state space is the location in the grid. This specific image was taken from [this link](#).

This domain is a 4×4 grid, where an agent starts at $(0, 0)$ and tries to navigate to $(3, 3)$. The states are the discrete positions in the grid and actions are the 4 cardinal directions. The reward function is 100 for reaching $(3, 3)$, -10 for reaching $(1, 1)$, 1 for reaching $(1, 3)$, and -1 for reaching all other states. If an agent takes an action that hits a wall, the agent stays in the same location. The transition dynamics are controlled by a parameter, p , where with probability p , an agent takes the intended action, else it takes an adjacent action with probability $(1 - p)/2$. All policies use a softmax action selection distribution with value θ_{sa} , for each state, s , and action a . The probability of taking action a in state s is given by:

$$\pi(a|s) = \frac{e^{\theta_{sa}}}{\sum_{a' \in \mathcal{A}} e^{\theta_{sa'}}$$

In the on-policy experiments, the evaluation and behavior policies were equiprobable policies in each cardinal direction. In the off-policy experiments, the evaluation policy was such that each θ was generated from a standard normal distribution and behavior policy was the equiprobable policy.

For the comparisons of batch linear PSEC-TD(0) and TD(0), we conducted a parameter sweep of the learning rates for the varying batch sizes. The parameter sweep was over: $\{5e^{-3}, 1e^{-3}, 5e^{-2}, 1e^{-2}, 5e^{-1}\}$. We used a value function convergence threshold of $1e^{-10}$. For PSEC-LSTD and LSTD, we stabilized the matrix, A , before inverting it by adding ϵI to the computed A . We conducted a parameter sweep over the following: $\epsilon \in \{1e^{-6}, 1e^{-5}, 1e^{-4}, 1e^{-3}, 1e^{-2}, 1e^{-1}\}$.

G.2. CartPole

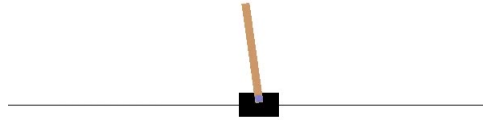


Figure 16. CartPole-v0 from OpenAI Gym (Brockman et al., 2016)

In this domain, the goal of the agent is to balance a pole for as long as possible. We trained our behavior policy using REINFORCE (Williams, 1992) with the Adam optimizer (Kingma & Ba, 2014) with learning rate $3e^{-4}$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. The behavior policy mapped raw state features to a softmax distribution over actions. The policy was a neural network with 2 hidden layers with 16 neurons each, and used the tanh activation function and was initialized with Xavier initialization (Glorot & Bengio, 2010).

The value function used by all algorithms was initialized by Xavier initialization and used the tanh activation function, and was trained using semi-gradient TD (Sutton & Barto, 2018). We used a convergence threshold of 0.1.

The PSEC policy was initialized by Xavier initialization and used the tanh activation function. PSEC-TD swept over the following learning rates $\alpha \in \{0.1 \times 2.0^j | j = -7, -6, \dots, 1, 2\}$. It used a validation set of 10% the size of the batch size. It used an L2 regularization of $2e^{-2}$. More details can be found in Section 5 and Appendix F.2.

G.3. InvertedPendulum

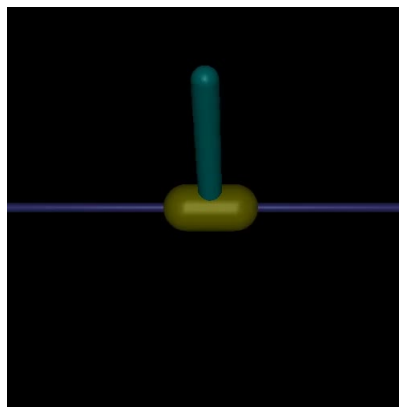


Figure 17. InvertedPendulum-v2 from OpenAI Gym and MuJoCo (Brockman et al., 2016; Todorov et al., 2012)

In this domain, the goal of the agent is to balance a pole for as long as possible. We trained our behavior policy using PPO (Schulman et al., 2017) with the default settings found on Gym (Brockman et al., 2016). The policy was a neural network with 2 hidden layers with 64 neurons each, and used the tanh activation function and was initialized with Xavier initialization (Glorot & Bengio, 2010). It mapped state features to an output vector that represented the mean vector of a Gaussian distribution. This mapping along with a separate parameter set representing the log standard deviation of each element in the output vector, make up the policy. The policy was trained by minimizing the following loss function:

$$\mathcal{L} = \sum_{i=1}^m 0.5((a_i - \mu(s_i))/e^\sigma)^2 + \sigma$$

where m are the number of state-action training examples, a_i is the action vector of the i^{th} example, $\mu(s_i)$ is the mean vector outputted by the neural network of the Gaussian distribution for state s_i , and σ is the the seperate parameter representing the log standard deviation of each element in the output vector, $\mu(s_i)$.

The value function used by all algorithms was initialized by Xavier initialization and used the tanh activation function, and was trained using semi-gradient TD (Sutton & Barto, 2018). We used a convergence threshold of 0.1.

The PSEC policy was initialized by Xavier initialization and used the tanh activation function. PSEC-TD swept over the following learning rates $\alpha \in \{0.1 \times 2.0^j | j = -8, -6, \dots, 1, 2\}$. It used a validation set of 20% the size of the batch size. More details can be found in Section 5 and Appendix F.2.