

---

# Policy Evaluation Using the $\Omega$ -Return

---

**Philip S. Thomas**

University of Massachusetts Amherst  
Carnegie Mellon University

**Scott Niekum**

University of Texas at Austin

**Georgios Theodorou**

Adobe Research

**George Konidaris**

Duke University

## Abstract

We propose the  $\Omega$ -return as an alternative to the  $\lambda$ -return currently used by the TD( $\lambda$ ) family of algorithms. The benefit of the  $\Omega$ -return is that it accounts for the correlation of different length returns. Because it is difficult to compute exactly, we suggest one way of approximating the  $\Omega$ -return. We provide empirical studies that suggest that it is superior to the  $\lambda$ -return and  $\gamma$ -return for a variety of problems.

## 1 Introduction

Most *reinforcement learning* (RL) algorithms learn a *value function*—a function that estimates the expected return obtained by following a given policy from a given state. Efficient algorithms for estimating the value function have therefore been a primary focus of RL research. The most widely used family of RL algorithms, the TD( $\lambda$ ) family [1], forms an estimate of return (called the  $\lambda$ -return) that blends low-variance but biased temporal difference return estimates with high-variance but unbiased Monte Carlo return estimates, using a parameter  $\lambda \in [0, 1]$ . While several different algorithms exist within the TD( $\lambda$ ) family—the original linear-time algorithm [1], least-squares formulations [2], and methods for adapting  $\lambda$  [3], among others—the  $\lambda$ -return formulation has remained unchanged since its introduction in 1988 [1].

Recently Konidaris et al. [4] proposed the  $\gamma$ -return as an alternative to the  $\lambda$ -return, which uses a more accurate model of how the variance of a return increases with its length. However, both the  $\gamma$  and  $\lambda$ -returns fail to account for the correlation of returns of different lengths, instead treating them as statistically independent. We propose the  $\Omega$ -return, which uses well-studied statistical techniques to directly account for the correlation of returns of different lengths. However, unlike the  $\lambda$  and  $\gamma$ -returns, the  $\Omega$ -return is not simple to compute, and often can only be approximated. We propose a method for approximating the  $\Omega$ -return, and show that it outperforms the  $\lambda$  and  $\gamma$ -returns on a range of off-policy evaluation problems.

## 2 Complex Backups

Estimates of return lie at the heart of value-function based RL algorithms: an estimate,  $\hat{V}^\pi$ , of the value function,  $V^\pi$ , estimates return from each state, and the learning process aims to reduce the error between estimated and observed returns. For brevity we suppress the dependencies of  $V^\pi$  and  $\hat{V}^\pi$  on  $\pi$  and write  $V$  and  $\hat{V}$ . *Temporal difference* (TD) algorithms use an estimate of the return obtained by taking a single transition in the *Markov decision process* (MDP) [5] and then estimating the remaining return using the estimate of the value function:

$$R_{s_t}^{\text{TD}} = r_t + \gamma \hat{V}(s_{t+1}),$$

where  $R_{s_t}^{\text{TD}}$  is the return estimate from state  $s_t$ ,  $r_t$  is the reward for going from  $s_t$  to  $s_{t+1}$  via action  $a_t$ , and  $\gamma \in [0, 1]$  is a discount parameter. Monte Carlo algorithms (for episodic tasks) do not use intermediate estimates but instead use the full return,

$$R_{s_t}^{\text{MC}} = \sum_{i=0}^{L-1} \gamma^i r_{t+i},$$

for an episode  $L$  transitions in length after time  $t$  (we assume that  $L$  is finite). These two types of return estimates can be considered instances of the more general notion of an  $n$ -step return,

$$R_{s_t}^{(n)} = \left( \sum_{i=0}^{n-1} \gamma^i r_{t+i} \right) + \gamma^n \hat{V}(s_{t+n}),$$

for  $n \geq 1$ . Here,  $n$  transitions are observed from the MDP and the remaining portion of return is estimated using the estimate of the value function. Since  $s_{t+L}$  is a state that occurs after the end of an episode, we assume that  $\hat{V}(s_{t+L}) = 0$ , always.

A *complex return* is a weighted average of the  $1, \dots, L$  step returns:

$$R_{s_t}^{\dagger} = \sum_{n=1}^L w_{\dagger}(n, L) R_{s_t}^{(n)}, \quad (1)$$

where  $w_{\dagger}(n, L)$  are weights and  $\dagger \in \{\lambda, \gamma, \Omega\}$  will be used to specify the weighting schemes of different approaches. The question that this paper proposes an answer to is: what weighting scheme will produce the best estimates of the true expected return?

The  $\lambda$ -return,  $R_{s_t}^{\lambda}$ , is the weighting scheme that is used by the entire family of  $\text{TD}(\lambda)$  algorithms [5]. It uses a parameter  $\lambda \in [0, 1]$  that determines how the weight given to a return decreases as the length of the return increases:

$$w_{\lambda}(n, L) = \begin{cases} (1 - \lambda)\lambda^{n-1} & \text{if } n < L \\ 1 - \sum_{i=1}^{n-1} w_{\lambda}(i) & \text{if } n = L. \end{cases}$$

When  $\lambda = 0$ ,  $R_{s_t}^{\lambda} = R_{s_t}^{\text{TD}}$ , which has low variance but high bias. When  $\lambda = 1$ ,  $R_{s_t}^{\lambda} = R_{s_t}^{\text{MC}}$ , which has high variance but is unbiased. Intermediate values of  $\lambda$  blend the high-bias but low-variance estimates from short returns with the low-bias but high-variance estimates from the longer returns.

The success of the  $\lambda$ -return is largely due to its simplicity— $\text{TD}(\lambda)$  using linear function approximation has per-time-step time complexity linear in the number of features. However, this efficiency comes at a cost: the  $\lambda$ -return is not founded on a principled statistical derivation.<sup>1</sup> Konidaris et al. [4] remedied this recently by showing that the  $\lambda$ -return is the maximum likelihood estimator of  $V(s_t)$  given three assumptions. Specifically,  $R_{s_t}^{\lambda} \in \arg \max_{x \in \mathbb{R}} \Pr(R_{s_t}^{(1)}, R_{s_t}^{(2)}, \dots, R_{s_t}^{(L)} | V(s_t) = x)$  if

**Assumption 1** (Independence).  $R_{s_t}^{(1)}, \dots, R_{s_t}^{(L)}$  are independent random variables,

**Assumption 2** (Unbiased Normal Estimators).  $R_{s_t}^{(n)}$  is normally distributed with mean  $\mathbf{E}[R_{s_t}^{(n)}] = V(s_t)$  for all  $n$ .

**Assumption 3** (Geometric Variance).  $\text{Var}(R_{s_t}^{(n)}) \propto 1/\lambda^n$ .

Although this result provides a theoretical foundation for the  $\lambda$ -return, it is based on three typically false assumptions: the returns are highly correlated, only the Monte Carlo return is unbiased, and the variance of the  $n$ -step returns from each state do not usually increase geometrically. This suggests three areas where the  $\lambda$ -return might be improved—it could be modified to better account for the correlation of returns, the bias of the different returns, and the true form of  $\text{Var}(R_{s_t}^{(n)})$ .

The  $\gamma$ -return uses an approximate formula for the variance of an  $n$ -step return in place of Assumption 3. This allows the  $\gamma$ -return to better account for how the variance of returns increases with their

<sup>1</sup>To be clear: there is a wealth of theoretical and empirical analyses of algorithms that use the  $\lambda$ -return. Until recently there was *not* a derivation of the  $\lambda$ -return as the estimator of  $V(s_t)$  that optimizes some objective (e.g., maximizes log likelihood or minimizes expected squared error).

length, while simultaneously removing the need for the  $\lambda$  parameter. The  $\gamma$ -return is given by the weighting scheme:

$$w_\gamma(n, L) = \frac{(\sum_{i=1}^n \gamma^{2(i-1)})^{-1}}{\sum_{\hat{n}=1}^L (\sum_{i=1}^{\hat{n}} \gamma^{2(i-1)})^{-1}}.$$

### 3 The $\Omega$ -Return

We propose a new complex return, the  $\Omega$ -return, that improves upon the  $\lambda$  and  $\gamma$  returns by accounting for the correlations of the returns. To emphasize this problem, notice that  $R_{s_t}^{(20)}$  and  $R_{s_t}^{(21)}$  will be almost identical (perfectly correlated) for many MDPs (particularly when  $\gamma$  is small). This means that Assumption 1 is particularly egregious, and suggests that a new complex return might improve upon the  $\lambda$  and  $\gamma$ -returns by properly accounting for the correlation of returns.

We formulate the problem of how best to combine different length returns to estimate the true expected return as a linear regression problem. This reformulation allows us to leverage the well-understood properties of linear regression algorithms. Consider a regression problem with  $L$  points,  $\{(x_i, y_i)\}_{i=1}^L$ , where the value of  $y_i$  depends on the value of  $x_i$ . The goal is to predict  $y_i$  given  $x_i$ . We set  $x_i = 1$  and  $y_i = R_{s_t}^{(i)}$ . We can then construct the design matrix (a vector in this case),  $\mathbf{x} = \mathbf{1} = [1, \dots, 1]^\top \in \mathbb{R}^L$  and the response vector,  $\mathbf{y} = [R_{s_t}^{(1)}, R_{s_t}^{(2)}, \dots, R_{s_t}^{(L)}]^\top$ . We seek a regression coefficient,  $\hat{\beta} \in \mathbb{R}$ , such that  $\mathbf{y} \approx \mathbf{x}\hat{\beta}$ . This  $\hat{\beta}$  will be our estimate of the true expected return.

*Generalized least squares* (GLS) is a method for selecting  $\hat{\beta}$  when the  $y_i$  are not necessarily independent and may have different variances. Specifically, if we use a linear model with (possibly correlated) mean-zero noise to model the data, i.e.,  $\mathbf{y} = \mathbf{x}\beta + \epsilon$ , where  $\beta \in \mathbb{R}$  is unknown,  $\epsilon$  is a random vector,  $\mathbb{E}[\epsilon] = \mathbf{0}$ , and  $\text{Var}(\epsilon|\mathbf{x}) = \Omega$ , then the GLS estimator

$$\hat{\beta} = (\mathbf{x}^\top \Omega^{-1} \mathbf{x})^{-1} \mathbf{x}^\top \Omega^{-1} \mathbf{y}, \quad (2)$$

is the *best linear unbiased estimator* (BLUE) for  $\beta$  [6]—the linear unbiased estimator with the lowest possible variance.

In our setting the assumptions about the true model that produced the data become that  $[R_{s_t}^{(1)}, R_{s_t}^{(2)}, \dots, R_{s_t}^{(L)}]^\top = [V(s_t), V(s_t), \dots, V(s_t)]^\top + \epsilon$ , where  $\mathbb{E}[\epsilon] = \mathbf{0}$  (i.e., the returns are all unbiased estimates of the true expected return) and  $\text{Var}(\epsilon|\mathbf{x}) = \Omega$ . Since  $\mathbf{x} = \mathbf{1}$  in our case,  $\text{Var}(\epsilon|\mathbf{x})(i, j) = \text{Cov}(R_{s_t}^{(i)} - V(s_t), R_{s_t}^{(j)} - V(s_t)) = \text{Cov}(R_{s_t}^{(i)}, R_{s_t}^{(j)})$ , where  $\text{Var}(\epsilon|\mathbf{x})(i, j)$  denotes the element of  $\text{Var}(\epsilon|\mathbf{x})$  in the  $i$ th row and  $j$ th column.

So, using only Assumption 2, GLS ((2), solved for  $\hat{\beta}$ ) gives us the complex return:

$$\hat{\beta} = \underbrace{\left( [1 \ 1 \ \dots \ 1] \Omega^{-1} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \right)^{-1}}_{=\frac{1}{\sum_{n,m=1}^L \Omega^{-1}(n,m)}} \underbrace{[1 \ 1 \ \dots \ 1] \Omega^{-1} \begin{bmatrix} R_{s_t}^{(1)} \\ R_{s_t}^{(2)} \\ \vdots \\ R_{s_t}^{(L)} \end{bmatrix}}_{=\sum_{n,m=1}^L \Omega^{-1}(n,m) R_{s_t}^{(n)}},$$

which can be written in the form of (1) with weights:

$$w_\Omega(n, L) = \frac{\sum_{m=1}^L \Omega^{-1}(n, m)}{\sum_{\hat{n}=1}^L \Omega^{-1}(\hat{n}, m)}, \quad (3)$$

where  $\Omega$  is an  $L \times L$  matrix with  $\Omega(i, j) = \text{Cov}(R_{s_t}^{(i)}, R_{s_t}^{(j)})$ .

Notice that the  $\Omega$ -return is a generalization of the  $\lambda$  and  $\gamma$  returns. The  $\lambda$ -return can be obtained by reintroducing the false assumption that the returns are independent and that their variance grows geometrically, i.e., by making  $\Omega$  a diagonal matrix with  $\Omega_{n,n} = \lambda^{-n}$ . Similarly, the  $\gamma$ -return can be obtained by making  $\Omega$  a diagonal matrix with  $\Omega_{n,n} = \sum_{i=1}^n \gamma^{2(i-1)}$ .

Notice that  $R_{s_t}^\Omega$  is a BLUE of  $V(s_t)$  if Assumption 2 holds. Since Assumption 2 does not hold, the  $\Omega$ -return is *not* an unbiased estimator of  $V(s)$ . Still, we expect it to outperform the  $\lambda$  and  $\gamma$ -returns because it accounts for the correlation of  $n$ -step returns and they do not. However, in some cases it may perform worse because it is still based on the false assumption that all of the returns are unbiased estimators of  $V(s_t)$ . Furthermore, given Assumption 2, there may be *biased* estimators of  $V(s_t)$  that have lower expected mean squared error than a BLUE (which must be unbiased).

## 4 Approximating the $\Omega$ -Return

In practice the covariance matrix,  $\Omega$ , is unknown and must be approximated from data. This approach, known as *feasible generalized least squares* (FGLS), can perform worse than ordinary least squares given insufficient data to accurately estimate  $\Omega$ . We must therefore accurately approximate  $\Omega$  from small amounts of data.

To study the accuracy of covariance matrix estimates, we estimated  $\Omega$  using a large number of trajectories for four different domains: a  $5 \times 5$  gridworld, a variant of the canonical mountain car domain, a real-world digital marketing problem, and a continuous control problem (DAS1), all of which are described in more detail in subsequent experiments. The covariance matrix estimates are depicted in Figures 1(a), 2(a), 3(a), and 4(a). We do not specify rows and columns in the figures because all covariance matrices and estimates thereof are symmetric. Because they were computed from a very large number of trajectories, we will treat them as ground truth.

We must estimate the  $\Omega$ -return when only a few trajectories are available. Figures 1(b), 2(b), 3(b), and 4(b) show direct empirical estimates of the covariance matrices using only a few trajectories. These empirical approximations are poor due to the very limited amount of data, except for the digital marketing domain, where a “few” trajectories means 10,000. The solid black entries in Figures 1(f), 2(f), 3(f), and 4(f) show the weights,  $w_\Omega(n, L)$ , on different length returns when using different estimates of  $\Omega$ . The noise in the direct empirical estimate of the covariance matrix using only a few trajectories leads to poor estimates of the return weights.

When approximating  $\Omega$  from a small number of trajectories, we must be careful to avoid this overfitting of the available data. One way to do this is to assume a compact parametric model for  $\Omega$ . Below we describe a parametric model of  $\Omega$  that has only four parameters, regardless of  $L$  (which determines the size of  $\Omega$ ). We use this parametric model in our experiments as a proof of concept—we show that the  $\Omega$ -return using even this simple estimate of  $\Omega$  can produce improved results over the other existing complex returns. We do not claim that this scheme for estimating  $\Omega$  is particularly principled or noteworthy.

### 4.1 Estimating Off-Diagonal Entries of $\Omega$

Notice in Figures 1(a), 2(a), 3(a), and 4(a) that for  $j > i$ ,  $\text{Cov}(R_{s_t}^i, R_{s_t}^j) \approx \text{Cov}(R_{s_t}^i, R_{s_t}^i) = \text{Var}(R_{s_t}^i)$ . This structure would mean that we can fill in  $\Omega$  given its diagonal values, leaving only  $L$  parameters. We now explain why this relationship is reasonable in general, and not just an artifact of our domains. We can write each entry in  $\Omega$  as a recurrence relation:

$$\begin{aligned} \text{Cov}[R_{s_t}^{(i)}, R_{s_t}^{(j)}] &= \text{Cov}[R_{s_t}^{(i)}, R_{s_t}^{(j-1)}] + \gamma^{j-1}(r_{t+j} + \gamma\hat{V}(s_{t+j}) - \hat{V}(s_{t+j-1})) \\ &= \text{Cov}[R_{s_t}^{(i)}, R_{s_t}^{(j-1)}] + \gamma^{j-1}\text{Cov}[R_{s_t}^{(i)}, r_{t+j} + \gamma\hat{V}(s_{t+j}) - \hat{V}(s_{t+j-1})], \end{aligned}$$

when  $i < j$ . The term  $r_{t+j} + \gamma\hat{V}(s_{t+j}) - \hat{V}(s_{t+j-1})$  is the temporal difference error  $j$  steps in the future. The proposed assumption that  $\text{Cov}(R_{s_t}^i, R_{s_t}^j) = \text{Var}(R_{s_t}^i)$  is equivalent to assuming that the covariance of this temporal difference error and the  $i$ -step return is negligible:  $\gamma^{j-1}\text{Cov}[R_{s_t}^{(i)}, r_{t+j} + \gamma\hat{V}(s_{t+j}) - \hat{V}(s_{t+j-1})] \approx 0$ . The approximate independence of these two terms is reasonable in general due to the Markov property, which ensures that at least the conditional covariance,  $\text{Cov}[R_{s_t}^{(i)}, r_{t+j} + \gamma\hat{V}(s_{t+j}) - \hat{V}(s_{t+j-1}) | s_t]$ , is zero.

Because this relationship is not exact, the off-diagonal entries tend to grow as they get farther from the diagonal. However, especially when some trajectories are padded with absorbing states, this relationship is quite accurate when  $j = L$ , since the temporal difference errors at the absorbing state are all zero, and  $\text{Cov}[R_{s_t}^{(i)}, 0] = 0$ . This results in a significant difference between  $\text{Cov}[R_{s_t}^{(i)}, R_{s_t}^{(L-1)}]$

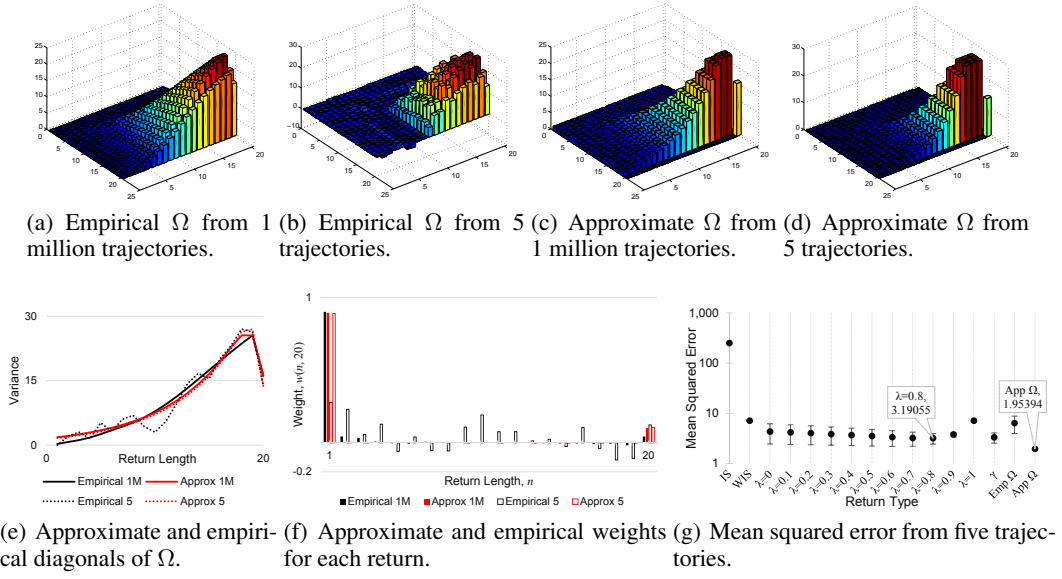


Figure 1: Gridworld Results.

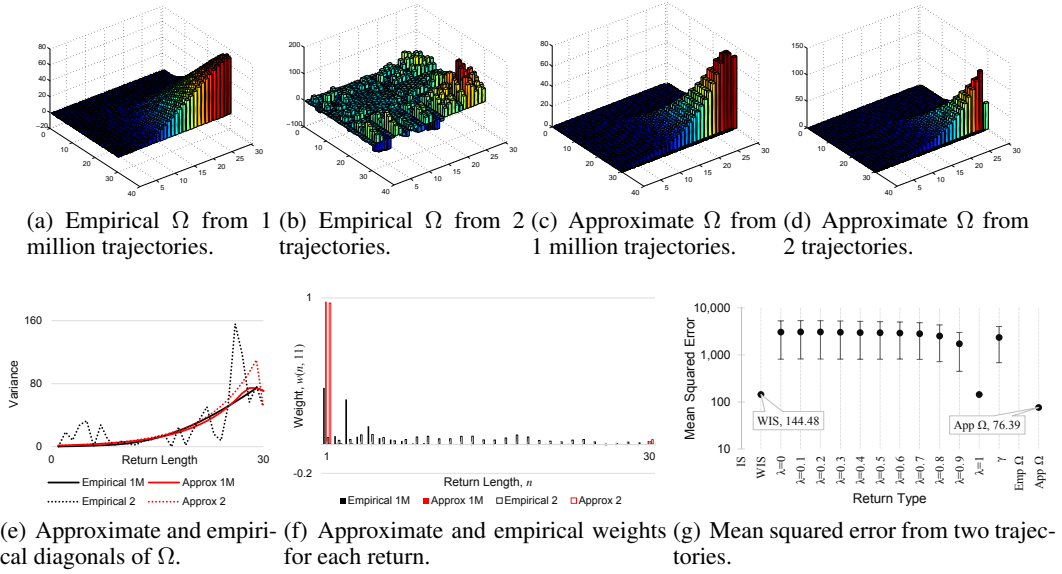


Figure 2: Mountain Car Results.

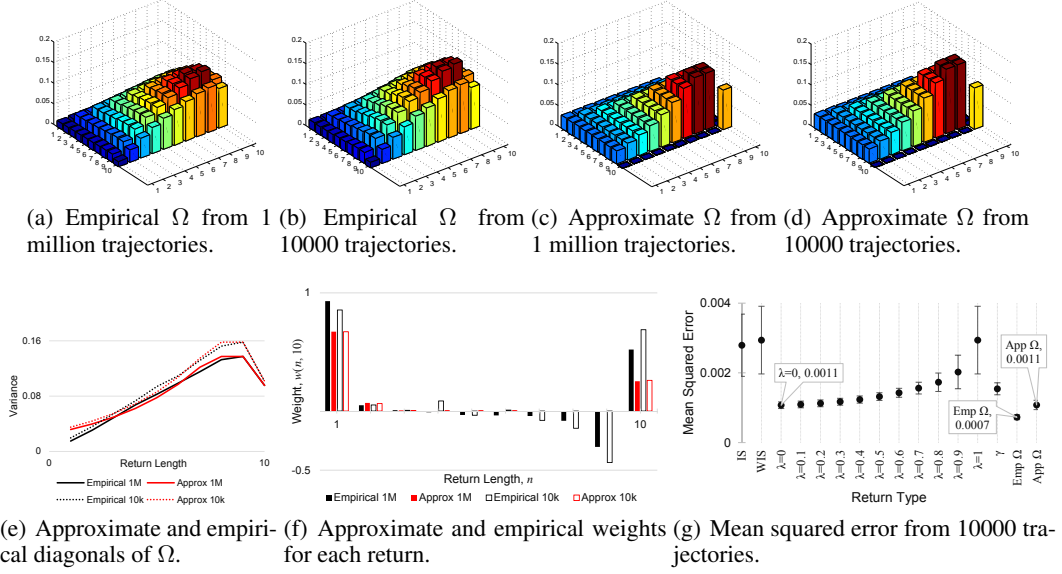


Figure 3: Digital Marketing Results.

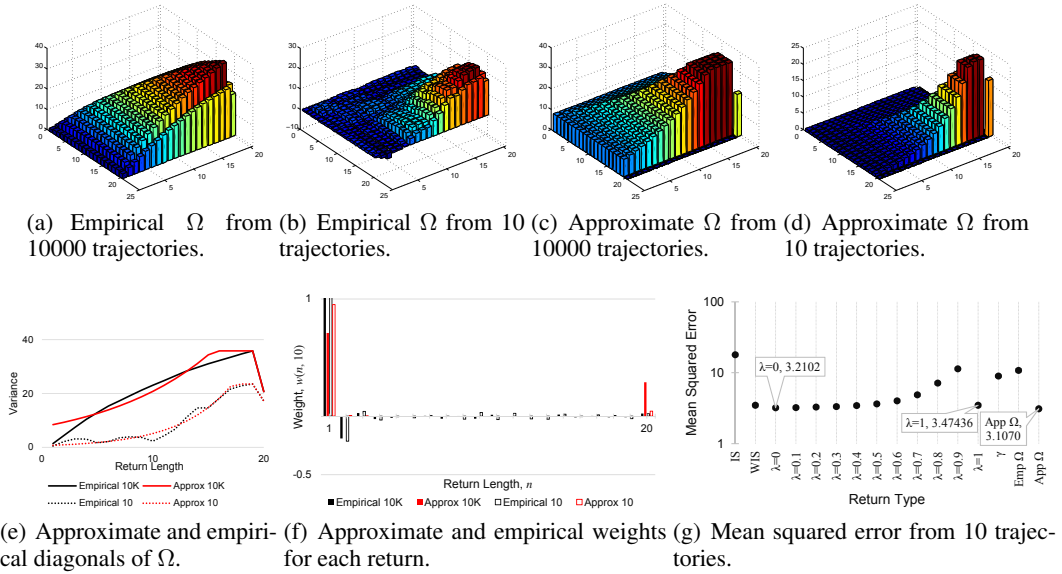


Figure 4: Functional Electrical Stimulation Results.

and  $\text{Cov}[R_{s_t}^{(i)}, R_{s_t}^{(L)}]$ . Rather than try to model this drop, which can influence the weights significantly, we reintroduce the assumption that the Monte Carlo return is independent of the other returns, making the off-diagonal elements of the last row and column zero.

## 4.2 Estimating Diagonal Entries of $\Omega$

The remaining question is how best to approximate the diagonal of  $\Omega$  from a very small number of trajectories. Consider the solid and dotted black curves in Figures 1(e), 2(e), 3(e), and 4(e), which depict the diagonals of  $\Omega$  when estimated from either a large number or small number of trajectories. When using only a few trajectories, the diagonal includes fluctuations that can have significant impacts on the resulting weights. However, when using many trajectories (which we treat as giving ground truth), the diagonal tends to be relatively smooth and monotonically increasing until it plateaus (ignoring the final entry).

This suggests using a smooth parametric form to approximate the diagonal, which we do as follows. Let  $v_i$  denote the sample variance of  $R_{s_t}^{(i)}$  for  $i = 1 \dots L$ . Let  $v_+$  be the largest sample variance:  $v_+ = \max_{i \in \{1, \dots, L\}} v_i$ . We parameterize the diagonal using four parameters,  $k_1$ ,  $k_2$ ,  $v_+$ , and  $v_L$ :

$$\hat{\Omega}_{k_1, k_2, v_+, v_L}(i, i) = \begin{cases} k_1 & \text{if } i = 1 \\ v_L & \text{if } i = L \\ \min\{v_+, k_1 k_2^{(1-i)}\} & \text{otherwise.} \end{cases}$$

$\Omega(1, 1) = k_1$  sets the initial variance, and  $v_L$  is the variance of the Monte Carlo return. The parameter  $v_+$  enforces a ceiling on the variance of the  $i$ -step return, and  $k_2$  captures the growth rate of the variance, much like  $\lambda$ . We select the  $k_1$  and  $k_2$  that minimize the mean squared error between  $\hat{\Omega}(i, i)$  and  $v_i$ , and set  $v_+$  and  $v_L$  directly from the data.<sup>2</sup>

This reduces the problem of estimating  $\Omega$ , an  $L \times L$  matrix, to estimating four numbers from return data. Consider Figures 1(c), 2(c), 3(c), and 4(c), which depict  $\hat{\Omega}$  as computed from many trajectories. The differences between these estimates and the ground truth show that this parameterization is not perfect, as we cannot represent the true  $\Omega$  exactly. However, the estimate is reasonable and the resulting weights (solid red) are visually similar to the ground truth weights (solid black) in Figures 1(f), 2(f), 3(f), and 4(f). We can now get accurate estimates of  $\Omega$  from very few trajectories. Figures 1(d), 2(d), 3(d), and 4(d) show  $\hat{\Omega}$  when computed from only a few trajectories. Note their similarity to  $\hat{\Omega}$  when using a large number of trajectories, and that the resulting weights (unfilled red in Figures 1(f), 2(f), 3(f), and 4(f)) are similar to the those obtained using many more trajectories (the filled red bars).

Pseudocode for approximating the  $\Omega$ -return is provided in Algorithm 1. Unlike the  $\lambda$ -return, which can be computed from a single trajectory, the  $\Omega$ -return requires a set of trajectories in order to estimate  $\Omega$ . The pseudocode assumes that every trajectory is of length  $L$ , which can be achieved by padding shorter trajectories with absorbing states.

<sup>2</sup>We include the constraints that  $k_2 \in [0, 1]$  and  $0 \leq k_1 \leq v_+$ .



---

**Algorithm 1:** Computing the  $\Omega$ -return.

---

**Require:**  $n$  trajectories beginning at  $s$  and of length  $L$ .

1. Compute  $R_s^{(i)}$  for  $i = 1, \dots, L$  and for each trajectory.
  2. Compute the sample variances,  $v_i = \text{Var}(R_s^{(i)})$ , for  $i = 1, \dots, L$ .
  3. Set  $v_+ = \max_{i \in \{1, \dots, L\}} v_i$ .
  4. Search for the  $k_1$  and  $k_2$  that minimize the mean squared error between  $v_i$  and  $\hat{\Omega}_{k_1, k_2, v_+, v_L}(i, i)$  for  $i = 1, \dots, L$ .
  5. Fill the diagonal of the  $L \times L$  matrix,  $\Omega$ , with  $\Omega(i, i) = \hat{\Omega}_{k_1, k_2, v_+, v_L}(i, i)$ , using the optimized  $k_1$  and  $k_2$ .
  6. Fill all of the other entries with  $\Omega(i, j) = \Omega(i, i)$  where  $j > i$ . If  $(i = L \text{ or } j = L)$  and  $i \neq j$  then set  $\Omega(i, j) = 0$  instead.
  7. Compute the weights for the returns according to (3).
  8. Compute the  $\Omega$ -return for each trajectory according to (1).
- 

## 5 Experiments

Approximations of the  $\Omega$ -return could, in principle, replace the  $\lambda$ -return in the whole family of  $\text{TD}(\lambda)$  algorithms. However, using the  $\Omega$ -return for  $\text{TD}(\lambda)$  raises several interesting questions that are beyond the scope of this initial work (e.g., is there a linear-time way to estimate the  $\Omega$ -return? Since a different  $\Omega$  is needed for every state, how can the  $\Omega$ -return be used with function approximation where most states will never be revisited?). We therefore focus on the specific problem of *off-policy policy evaluation*—estimating the performance of a policy using trajectories generated by a possibly different policy. This problem is of interest for applications that require the evaluation of a proposed policy using historical data.

Due to space constraints, we relegate the details of our experiments to the appendix in the supplemental documents. However, the results of the experiments are clear—Figures 1(g), 2(g), 3(g), and 4(g) show the *mean squared error* (MSE) of value estimates when using various methods.<sup>3</sup> Notice that, for all domains, using the  $\Omega$ -return (the EMP  $\Omega$  and APP  $\Omega$  labels) results in lower MSE than the  $\gamma$ -return and the  $\lambda$ -return with any setting of  $\lambda$ .

## 6 Conclusions

Recent work has begun to explore the statistical basis of complex estimates of return, and how we might reformulate them to be more statistically efficient [4]. We have proposed a return estimator that improves upon the  $\lambda$  and  $\gamma$ -returns by accounting for the covariance of return estimates. Our results show that understanding and exploiting the fact that in control settings—unlike in standard supervised learning—observed samples are typically neither independent nor identically distributed, can substantially improve data efficiency in an algorithm of significant practical importance.

Many (largely positive) theoretical properties of the  $\lambda$ -return and  $\text{TD}(\lambda)$  have been discovered over the past few decades. This line of research into other complex returns is still in its infancy, and so there are many open questions. For example, can the  $\Omega$ -return be improved upon by removing Assumption 2 or by keeping Assumption 2 but using a biased estimator (not a BLUE)? Is there a method for approximating the  $\Omega$ -return that allows for value function approximation with the same time complexity as  $\text{TD}(\lambda)$ , or which better leverages our knowledge that the environment is Markovian? Would  $\text{TD}(\lambda)$  using the  $\Omega$ -return be convergent in the same settings as  $\text{TD}(\lambda)$ ? While we hope to answer these questions in future work, it is also our hope that this work will inspire other researchers to revisit the problem of constructing a statistically principled complex return.

---

<sup>3</sup>To compute the MSE we used a large number of Monte Carlo rollouts to estimate the true value of each policy.



## References

- [1] R.S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1):9–44, 1988.
- [2] S.J. Bradtke and A.G. Barto. Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22(1-3):33–57, March 1996.
- [3] C. Downey and S. Sanner. Temporal difference Bayesian model averaging: A Bayesian perspective on adapting lambda. In *Proceedings of the 27th International Conference on Machine Learning*, pages 311–318, 2010.
- [4] G.D. Konidaris, S. Niekum, and P.S. Thomas.  $TD_{\gamma}$ : Re-evaluating complex backups in temporal difference learning. In *Advances in Neural Information Processing Systems 24*, pages 2402–2410, 2011.
- [5] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.
- [6] T. Kariya and H. Kurata. *Generalized Least Squares*. Wiley, 2004.
- [7] D. Precup, R. S. Sutton, and S. Singh. Eligibility traces for off-policy policy evaluation. In *Proceedings of the 17th International Conference on Machine Learning*, pages 759–766, 2000.
- [8] A. R. Mahmood, H. Hasselt, and R. S. Sutton. Weighted importance sampling for off-policy learning with linear function approximation. In *Advances in Neural Information Processing Systems 27*, 2014.
- [9] J. R. Tetreault and D. J. Litman. Comparing the utility of state features in spoken dialogue using reinforcement learning. In *Proceedings of the Human Language Technology/North American Association for Computational Linguistics*, 2006.
- [10] G. D. Konidaris, S. Osentoski, and P. S. Thomas. Value function approximation in reinforcement learning using the Fourier basis. In *Proceedings of the Twenty-Fifth Conference on Artificial Intelligence*, pages 380–395, 2011.
- [11] G. Theodorou and A. Hallak. Lifetime value marketing using reinforcement learning. In *The 1st Multidisciplinary Conference on Reinforcement Learning and Decision Making*, 2013.
- [12] P. S. Thomas, G. Theodorou, and M. Ghavamzadeh. High confidence off-policy evaluation. In *Proceedings of the Twenty-Ninth Conference on Artificial Intelligence*, 2015.
- [13] D. Blana, R. F. Kirsch, and E. K. Chadwick. Combined feedforward and feedback control of a redundant, nonlinear, dynamic musculoskeletal system. *Medical and Biological Engineering and Computing*, 47: 533–542, 2009.
- [14] P. S. Thomas, M. S. Branicky, A. J. van den Bogert, and K. M. Jagodnik. Application of the actor-critic architecture to functional electrical stimulation control of a human arm. In *Proceedings of the Twenty-First Innovative Applications of Artificial Intelligence*, pages 165–172, 2009.
- [15] P. M. Pilariski, M. R. Dawson, T. Degris, F. Fahimi, J. P. Carey, and R. S. Sutton. Online human training of a myoelectric prosthesis controller via actor-critic reinforcement learning. In *Proceedings of the 2011 IEEE International Conference on Rehabilitation Robotics*, pages 134–140, 2011.
- [16] K. Jagodnik and A. van den Bogert. A proportional derivative FES controller for planar arm movement. In *12th Annual Conference International FES Society*, Philadelphia, PA, 2007.
- [17] N. Hansen and A. Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9(2):159–195, 2001.

## Appendix 1: Off-Policy Evaluation

The performance of a policy is the expected return of a trajectory generated when following that policy, i.e., the value of the initial state.<sup>4</sup> Given a set  $D$  of  $m$  trajectories,  $\tau_i, i = 1, \dots, m$ , generated by a *behavior policy*,  $\pi_b$ , and an *evaluation policy*,  $\pi_e$ , we must estimate the expected return of  $\pi_e$ :  $V^{\pi_e}(s_0) = \mathbb{E}[R_{s_0}^{MC} | \pi_e]$ . We now describe existing methods for off-policy evaluation and then propose a new approach that uses the  $\Omega$ -return.

### Importance Sampling

If the trajectories were on-policy (i.e.,  $\pi_b = \pi_e$ ) we could estimate the expected return of  $\pi_e$  by averaging the observed returns. However, when the evaluation policy differs from the behavior policy, we must take a weighted average of the returns where larger weights are given to trajectories that are more likely under the evaluation policy. *Importance sampling* (IS), a principled way to perform this reweighting, produces an unbiased estimate of  $V^{\pi_e}(s_0)$  from each trajectory as:

$$\hat{V}_{MC}^{\pi_e}(s_0, \tau) = R_{s_0}^{MC} \underbrace{\prod_{t=0}^{L-1} \frac{\pi_e(s_t^\tau, a_t^\tau)}{\pi_b(s_t^\tau, a_t^\tau)}}_{\text{Importance Weight}},$$

where  $s_t^\tau$  and  $a_t^\tau$  are the state and action at time  $t$  in trajectory  $\tau$ ,  $L$  is the length of the longest trajectory,<sup>5</sup> and  $\pi(s, a)$  denotes the probability of action  $a$  in state  $s$  under policy  $\pi$ .  $\hat{V}_{MC}^{\pi_e}(s_0, \tau_i)$  is called the *importance weighted return*. The estimate of  $V^{\pi_e}(s_0)$  from  $D$  can then be produced by averaging the importance weighted return from each trajectory:

$$\hat{V}_{IS}^{\pi_e}(s_0, D) = \frac{1}{m} \sum_{\tau \in D} \hat{V}_{MC}^{\pi_e}(s_0, \tau).$$

Although  $\hat{V}_{IS}^{\pi_e}(s_0, D)$  is an unbiased estimate of  $V^{\pi_e}(s_0)$ , it has high variance since the importance weights can be large for long trajectories [7].

### Weighted Importance Sampling

In *weighted importance sampling* (WIS) the sum of the importance weighted returns is divided by the sum of importance weights rather than the number of trajectories:

$$\hat{V}_{WIS}^{\pi_e}(s_0, D) = \frac{1}{\sum_{\tau \in D} \prod_{t=0}^{L-1} \frac{\pi_e(s_t^\tau, a_t^\tau)}{\pi_b(s_t^\tau, a_t^\tau)}} \sum_{i=1}^n \hat{V}_{MC}^{\pi_e}(s_0, \tau_i).$$

This modification means that  $\hat{V}_{WIS}^{\pi_e}(s_0, D)$  is a *biased* but consistent estimator of  $V^{\pi_e}(s_0)$ . It also tends to have significantly lower variance than  $\hat{V}_{IS}^{\pi_e}(s_0, D)$ ; WIS trades bias for lower variance relative to IS, and is the current method of choice for off-policy evaluation [7, 8].

### Complex Weighted Importance Sampling

We propose a simple modification to WIS, which we call *Complex Weighted Importance Sampling* (CWIS). Whereas IS and WIS use the Monte Carlo return in  $\hat{V}_{MC}^{\pi_e}(s_0, \tau)$ , CWIS uses a complex return like  $R_{s_0}^\lambda$ ,  $R_{s_0}^\gamma$ , or  $R_{s_0}^\Omega$ . This replaces the unbiased Monte Carlo estimate of return with a lower variance but biased (and not necessarily consistent) estimator of expected return. This is a worthwhile trade-off because the primary drawback of IS and WIS is their high variance. We choose to focus on using complex returns for WIS rather than IS because the use of complex backups would make IS biased, damaging its primary selling point.

The CWIS<sub>†</sub> estimator of  $V^{\pi_e}(s_0)$  is:

$$\hat{V}_{CWIS_\dagger}^{\pi_e}(s_0, D) = \frac{1}{\sum_{\tau \in D} \prod_{t=0}^{L-1} \frac{\pi_e(s_t^\tau, a_t^\tau)}{\pi_b(s_t^\tau, a_t^\tau)}} \sum_{i=1}^n \hat{V}_\dagger^{\pi_e}(s_0, \tau_i), \quad (4)$$

where  $\hat{V}_\dagger^{\pi_e}(s_0, \tau) = R_{s_0}^\dagger \prod_{t=0}^{L-1} \frac{\pi_e(s_t^\tau, a_t^\tau)}{\pi_b(s_t^\tau, a_t^\tau)}$ , and  $\dagger$  specifies the complex return, e.g.,  $\lambda$ ,  $\gamma$ ,  $\Omega$ , or 0. Notice that CWIS<sub>MC</sub> is equivalent to WIS and CWIS <sub>$\lambda=1$</sub> . Also, CWIS <sub>$\lambda=0$</sub>  is the ECR estimator [9].

<sup>4</sup>A distribution over initial states can be handled by inserting a new initial state with only one admissible action, no reward, which transitions according to the start distribution.

<sup>5</sup>All shorter trajectories are padded to be of length  $L$  using a zero-reward absorbing state with one admissible action.

CWIS assumes that an estimate of the off-policy value function is available when computing returns of different lengths. Any off-policy evaluation method can be used to produce an estimate of the value function for  $\pi_e$ . Ideally this off-policy value function approximation would be computed to minimize the mean squared error between the  $\Omega$ -return and the approximation of the value function. However, value function approximation methods for the  $\Omega$ -return are not yet available, so in their absence we use *WIS-LSTD*(0.5) [8] with the Fourier basis [10].

An implementation of  $\text{CWIS}_\dagger$  has four steps. **1)** Estimate the off-policy value function. Here we used *WIS-LSTD*(0.5). **2)** Compute the different length returns  $R_{s_0}^{(n)}$  for each trajectory for  $n = 1, \dots, L$ . **3)** Compute the  $\dagger$ -return,  $R_{s_0}^\dagger$ , for each trajectory. **4)** Use (4) to compute the final estimate of the performance of  $\pi_e$ .

Note that if the complex return used by  $\text{CWIS}_\dagger$  is the same as that used by the off-policy value function approximation method, then  $\text{CWIS}_\dagger$  returns the *target* of the off-policy value function approximation method for start state  $s_0$ .

## Appendix 2: Off-Policy Evaluation Experiments

We compared IS, WIS,  $\text{CWIS}_0$ ,  $\text{CWIS}_\lambda$  (for various  $\lambda$ ),  $\text{CWIS}_\gamma$ , and  $\text{CWIS}_\Omega$  on four benchmark domains. In three cases  $\text{CWIS}_\Omega$  using our approximation to  $\Omega$  outperforms all other methods, and in the fourth case it performs as well as every other method tested except for  $\text{CWIS}_\Omega$  using a direct empirical estimate of  $\Omega$ .

**5 × 5 Gridworld.** Our first empirical study used a 5 × 5 gridworld [5] with a reasonable hand-tuned behavior policy and an improved evaluation policy. We sampled 1 million trajectories from the evaluation policy and used the average return as the ground truth value of the evaluation policy. We then repeated the following process 10000 times: sample 5 trajectories and use IS, WIS,  $\text{CWIS}_0$ ,  $\text{CWIS}_\lambda$  (with 11 different values of  $\lambda$ ),  $\text{CWIS}_\gamma$ , and  $\text{CWIS}_\Omega$  to estimate the value of the evaluation policy. The *mean squared error* (MSE) over these 10000 trials, including standard error bars, is shown in Figure 1(g). In this and subsequent similar figures, “Emp  $\Omega$ ” (empirical) denotes  $\text{CWIS}_\Omega$  using the sample covariance matrix to estimate  $\Omega$  while “App  $\Omega$ ” (approximate) denotes  $\text{CWIS}_\Omega$  using our proposed scheme for smoothly approximating  $\Omega$ . Notice that our scheme for approximating  $\Omega$  results in  $\text{CWIS}_\Omega$  having the lowest MSE, with  $\text{CWIS}_\lambda$  with  $\lambda = 0.8$  having the second lowest MSE.

**Mountain Car.** Our second study was set up similarly, using a simplified variant of the canonical mountain car domain [5], where each action was held constant for twenty time steps to shorten trajectories and avoid numerical overflows with IS. The behavior policy was random, the evaluation policy near-optimal, and only two trajectories were used for each estimate of the evaluation policy’s performance (this makes the problem challenging for CWIS, since the estimate of the value function is poor). The MSEs from this experiment are depicted in Figure 2(g).

Notice the large MSE, which are due to the disparity between the behavior and evaluation policies. This caused large importance weights which can cause instability in value function approximation methods like *WIS-LSTD* when using only two trajectories to estimate the value function. As a result, the accuracy of the approximate value function produced by *WIS-LSTD* can vary.  $\text{CWIS}_{\lambda=1}$ , which is equivalent to WIS, therefore performs well because it does not use the approximate value function produced by *WIS-LSTD*. However, “Approx  $\Omega$ ” performs best, selectively using the value function only when it is accurate enough to produce low variance returns.

**Digital Marketing.** Our third case study involved targeting advertisements. When a user visits a webpage, a vector of features describing the user is used to select an advertisement to display. A well-targeted advertisement will be of interest to the user and has a higher probability of being clicked than a poorly-targeted one. It has been shown that the problem of which advertisement to show should be treated as a sequential decision problem (as opposed to a bandit problem) [11].

We used the same digital marketing simulator as Thomas et al. [12], which was trained using real data from a Fortune 50 company, to evaluate a newly optimized policy using 10000 trajectories from a previous (slightly worse) policy. This process was repeated 100 times to produce Figure 3(g). Notice that here we are using significantly more trajectories to form each estimate of the evaluation policy’s performance. This is because the simulator was built using real data from a company that receives hundreds of thousands of unique visitors per day, so that 10000 trajectories is a “small” amount of data. With so many trajectories, we can use the sample covariance matrix to accurately estimate  $\Omega$  for  $\text{CWIS}_\Omega$ , and hence this approach (“Emp  $\Omega$ ”) performs best, though our approximate  $\Omega$  performs as well as the best setting of  $\lambda$ .

**Functional Electrical Stimulation.** Our fourth study used a simulation, *Dynamic Arm Simulator 1* [13], of a human arm undergoing *functional electrical stimulation* (FES). The goal of FES is to use direct stimulation of the muscles in a paralyzed subject’s arm to move the arm from the current position to a desired position. Researchers have used RL previously for simulated studies into the efficacy of RL for automatically optimiz-

ing the controller for each individual’s arm and to mitigate sources of nonstationarity including fatigue [14]. Researchers have also considered using RL for control of ordinary prosthetics [15].

Here we use the same arm model and domain specifications as Thomas et al. [14], except with longer time steps in order to decrease trajectory lengths to avoid numerical instabilities in IS. The behavior policy performs similar to a *proportional derivative* (PD) controller optimized for FES control of a human arm [16], and the evaluation policy is a slightly improved policy found using the CMA-ES policy search algorithm [17].

In a real deployment of FES, each trajectory is a two-second reaching movement, and so a controller that can adapt using only tens of trajectories is important. We therefore consider the problem of evaluating a proposed policy improvement using only ten trajectories. This is repeated 1000 times, and the resulting MSEs are provided in Figure 4(g).  $CWIS_{\Omega}$  using our approximation performs best.