# Policy Evaluation Using the ?-Return

**Authored by:**

Georgios Theocharous
George Konidaris
Scott Niekum
Philip S. Thomas

**Abstract**

We propose the ?-return as an alternative to the ?-return currently used by the TD(?) family of algorithms. The benefit of the ?-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 ?-return. We provide empirical studies that suggest that it is superior to the ?-return and ?-return for a variety of problems.

## 1 Paper Body

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(?) family [1], forms an estimate of return (called the ?-return) that blends low-variance but biased temporal difference return estimates with high-variance but unbiased Monte Carlo return estimates, using a parameter ? ? [0, 1]. While several different algorithms exist within the TD(?) family?the original linear-time algorithm [1], least-squares formulations [2], and methods for adapting ? [3], among others?the ?-return formulation has remained unchanged since its introduction in 1988 [1]. Recently Konidaris et al. [4] proposed the ?-return as an alternative to the ?-return, which uses a more accurate model of how the variance of a return increases with its length. However, both the ? and ?-returns fail to account for the correlation of returns of different lengths, instead treating them as statistically independent. We propose the ?-return, which uses well-studied statistical techniques to directly account for the correlation of returns of different lengths. However, unlike the ? and ?-returns, the ?-return is not simple to compute, and often can only be approximated. We propose a method for approximating the ?-return, and show that it outperforms the ? and ?-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, V? ? , of the value function, V ? , 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 ? and V? ? on ? and write V and 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: RsTD = rt + ? V? (st+1 ), t 1

is the return estimate from state st , rt is the reward for going from st to st+1 via action where RsTD t at , and ? ? [0, 1] is a discount parameter. Monte Carlo algorithms (for episodic tasks) do not use intermediate estimates but instead use the full return, RsMC = t

L?1 X
? i rt+i ,
i=0

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, ! n?1 X (n) i Rs = ? rt+i + ? n V? (st+n ), t

i=0

for n ? 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 st+L is a state that occurs after the end of an episode, we assume that V? (st+L ) = 0, always. A complex return is a weighted average of the 1, . . . , L step returns: Rs?t =

L X
w? (n, L)Rs(n) , t
(1)
n=1

where w? (n, L) are weights and ? ? {?, ?, ?} 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 ?-return, Rs?t , is the weighting scheme that is used by the entire family of TD(?) algorithms [5]. It uses a parameter ? ? [0, 1] that determines how the weight given to a return decreases as the length of the return increases: ( (1 ? ?)?n?1 if n ¡ L Pn?1 w? (n, L) = 1 ? i=1 w? (i) if n = L. , which , which has low variance but high bias. When ? = 1, Rs?t = RsMC When ? = 0, Rs?t = RsTD t t has high variance but is unbiased. Intermediate values of ? 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 ?-return is largely due to its simplicity?TD(?) using linear function approximation has per-time-step time complexity linear in the number of features. However, this efficiency comes at a cost: the ?-return is not founded on a principled statistical derivation.1 Konidaris et al. [4] remedied this recently by showing that the ?-return is the maximum likelihood estimator

of V (st ) (1) (2) (L) given three assumptions. Specifically, Rs?t ? arg maxx?R Pr(Rst , Rst , . . . , Rst —V (st ) = x) if (1)

(L)

Assumption 1 (Independence). Rst , . . . , Rst are independent random variables, (n)

(n)

Assumption 2 (Unbiased Normal Estimators). Rst is normally distributed with mean E[Rst ] = V (st ) for all n. (n)

Assumption 3 (Geometric Variance). Var(Rst ) ? 1/?n . Although this result provides a theoretical foundation for the ?-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 ?-return might be improved?it could be modified to better account for the (n) correlation of returns, the bias of the different returns, and the true form of Var(Rst ). The ?-return uses an approximate formula for the variance of an n-step return in place of Assumption 3. This allows the ?-return to better account for how the variance of returns increases with their 1 To be clear: there is a wealth of theoretical and empirical analyses of algorithms that use the ?-return. Until recently there was not a derivation of the ?-return as the estimator of V (st ) that optimizes some objective (e.g., maximizes log likelihood or minimizes expected squared error).

2

length, while simultaneously removing the need for the ? parameter. The ?-return is given by the weighting scheme: Pn ( i=1 ? 2(i?1) )?1 w? (n, L) = PL Pn? . 2(i?1) )?1 n ? =1 ( i=1 ?

3

The ?-Return

We propose a new complex return, the ?-return, that improves upon the ? and ? returns by account(20) (21) ing for the correlations of the returns. To emphasize this problem, notice that Rst and Rst will be almost identical (perfectly correlated) for many MDPs (particularly when ? is small). This means that Assumption 1 is particularly egregious, and suggests that a new complex return might improve upon the ? and ?-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 wellunderstood properties of linear regression algorithms. Consider a regression problem with L points, {(xi , yi )}L i=1 , where the value of yi depends on the value of xi . The goal is to predict yi given (i) xi . We set xi = 1 and yi = Rst . We can then construct the design matrix (a vector in this case), (1) (2) (L) — L x = 1 = [1, . . . , 1] ? R and the response vector, y = [Rst , Rst , . . . , Rst ]— . We seek a re? This ?? will be our estimate of the true expected gression coefficient, ?? ? R, such that y ? x?. return. Generalized least squares (GLS) is a method for selecting ?? when the yi are not necessarily independent and may have different variances. Specifically, if we use a linear model with (possibly correlated) mean-zero noise

3

to model the data, i.e., y = x? + , where ? ? R is unknown, is a random vector, E[] = 0, and Var(—x) = ?, then the GLS estimator ?? = (x— ??1 x)?1 x— ??1 y,

(2)

is the best linear unbiased estimator (BLUE) for ? [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 (1) (2) (L) [Rst , Rst , . . . , Rst ]— = [V (st ), V (st ), . . . , V (st )]— + , where E[] = 0 (i.e., the returns are all unbiased estimates of the true expected return) and Var(—x) = ?. Since x = 1 in our case, (i) (j) (i) (j) Var(—x)(i, j) = Cov(Rst ? V (st ), Rst ? V (st )) = Cov(Rst , Rst ), where Var(—x)(i, j) denotes the element of Var(—x) in the ith row and jth column. ? gives us the complex return: So, using only Assumption 2, GLS ((2), solved for ?) ? (1) ? ? ???1 Rst 1 ? (2) ? ? 1 ?? ? Rst ?1 ? ?1 ? ?? ?? = ? ?[ 1 1 . . . 1 ] ? ? ... ?? [ 1 1 . . . 1 ] ? ? ? .. ? . (L) 1 Rst — {z }— {z 1 = PL

?1 (n,m) n,m=1 ?

P ?1 (n,m)R(n) = L st n,m=1 ?

which can be written in the form of (1) with weights: PL ??1 (n, m) w? (n, L) = PLm=1 , ?1 (? n, m) n ? ,m=1 ? (i)

? ? ? ?, ? ? }

(3)

(j)

where ? is an L ? L matrix with ?(i, j) = Cov(Rst , Rst ). Notice that the ?-return is a generalization of the ? and ? returns. The ?-return can be obtained by reintroducing the false assumption that the returns are independent and that their variance grows geometrically, i.e., by making ? a diagonal matrix with ?n,n = ??n . Similarly, the ?-return can be Pn obtained by making ? a diagonal matrix with ?n,n = i=1 ? 2(i?1) . 3

Notice that Rs?t is a BLUE of V (st ) if Assumption 2 holds. Since Assumption 2 does not hold, the ?-return is not an unbiased estimator of V (s). Still, we expect it to outperform the ? and ?-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 (st ). Furthermore, given Assumption 2, there may be biased estimators of V (st ) that have lower expected mean squared error than a BLUE (which must be unbiased).

4

Approximating the ?-Return

In practice the covariance matrix, ?, 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 ?. We must therefore accurately approximate ? from small amounts of data. To study the accuracy of covariance matrix estimates, we estimated ? using a large number of trajectories for four different domains: a 5 ? 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 ?-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? (n, L), on different length returns when using different estimates of ?. 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 ? 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 ?. Below we describe a parametric model of ? that has only four parameters, regardless of L (which determines the size of ?). We use this parametric model in our experiments as a proof of concept? we show that the ?-return using even this simple estimate of ? can produce improved results over the other existing complex returns. We do not claim that this scheme for estimating ? is particularly principled or noteworthy.

4.1

Estimating Off-Diagonal Entries of ?

Notice in Figures 1(a), 2(a), 3(a), and 4(a) that for j ¿ i, Cov(Rsi t , Rsj t ) ? Cov(Rsi t , Rsi t ) = Var(Rsi t ). This structure would mean that we can fill in ? 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 ? as a recurrence relation: Cov[Rs(i) , Rs(j) ] =Cov[Rs(i) , Rs(j?1) + ? j?1 (rt+j + ? V? (st+j ) ? V? (st+j?1 )] t t t t =Cov[Rs(i) , Rs(j?1) ] + ? j?1 Cov[Rs(i) , rt+j + ? V? (st+j ) ? V? (st+j?1 )], t t t

when i ¡ j. The term rt+j + ? V? (st+j ) ? V? (st+j?1 ) is the temporal difference error j steps in the future. The proposed assumption that Cov(Rsi t , Rsj t ) = Var(Rsi t ) is equivalent to assuming that the covariance of this temporal difference error and the i-step return is negligible: (i) ? j?1 Cov[Rst , rt+j + ? V? (st+j ) ? V? (st+j?1 )] ? 0. The approximate independence of these two terms is reasonable in general due to the Markov property, which ensures that at least the conditional (i) covariance, Cov[Rst , rt+j + ? V? (st+j ) ? V? (st+j?1 )—st ], 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 (i) (i) (L?1) are all zero, and Cov[Rst , 0] = 0. This results in a significant difference between Cov[Rst , Rst ] 4

25

30

20
20
25
30
20 20
15
15 10
10
10 0
5
?10 0
0 0
0 0 5
5
20
10
20
10
15
15
10 20
5 25
25
5
15
15
10 20
5
20
10
15
15
10 20
0 0
5
20
10
15
15
10
5
10 20
5
5 25
25

(a) Empirical ? from 1 (b) Empirical ? from 5 (c) Approximate ? from (d) Approximate ? from million trajectories. trajectories. 1 million trajectories. 5 trajectories. 1

1,000
Mean Squared Error
Weight, w(n, 20)
Variance
30
15
0 0
1
20
Return Length Empirical 1M
Approx 1M
Empirical 5
Approx 5
100
20
-0.2
App ?, 1.95394
?=0.8, 3.19055
10
1
Return Length, n Empirical 1M
Approx 1M
Empirical 5
Return Type
Approx 5

(e) Approximate and empiri- (f) Approximate and empirical weights (g) Mean squared error from five trajeccal diagonals of ?. for each return. tories.

Figure 1: Gridworld Results.

80
200
80 60
150
60 100
100
40
40
20
0
0
?100 0
?20 0 10
0 0 10
30 20 15

7

30 5
40
30 20
5
40
15
30
10
5
25
20
20 15
30
10
10
30 25
20
20 15
30
10 40
25
20
0 0
10
30
25
20
50
20
10 5
40

(a) Empirical ? from 1 (b) Empirical ? from 2 (c) Approximate ? from (d) Approximate ? from million trajectories. trajectories. 1 million trajectories. 2 trajectories. 1

10,000
Approx 1M
Empirical 2
Approx 2
?
App ?
?=1
?=0.9
Emp ?
Return Type
?=0.8
?=0.7

?=0.6

Return Length, n Empirical 1M

?=0.5

-0.2

App ?, 76.39 ?=0.4

Approx 2

WIS, 144.48 10 ?=0.3

Approx 1M

Empirical 2

30

?=0.2

Return Length Empirical 1M

1

30

?=0

0

100

?=0.1

0

1,000

IS

80

WIS

Mean Squared Error

Weight, w(n, 11)

Variance

160

(e) Approximate and empiri- (f) Approximate and empirical weights (g) Mean squared error from two trajeccal diagonals of ?. for each return. tories.

Figure 2: Mountain Car Results.

5

0.2

0.2

0.2

0.2

0.15

0.15

0.15

0.15

0.1

0.1

0.1

0.1

0.05

0.05

0.05

0.05
0
0 1
2
4
5
6
7
8
9 10 1
2
3
4
5
6
7
8
9
0
0
1 3
2
10
1
3
4
5
6
7
8
9 10 1
3
2
4
5
6
8
7
9
2
10
1
3
4
5
6

7
8
9 10 1
2
3
4
7
6
5
9
8
2
3
10
4
5
6
7
8
9 10 1
2
3
4
5
6
8
7
10
9

(a) Empirical ? from 1 (b) Empirical ? from (c) Approximate ? from (d) Approximate ? from million trajectories. 10000 trajectories. 1 million trajectories. 10000 trajectories. 1

0.08
1
Mean Squared Error
0.004
Weight, w(n, 10)
10
App ?, 0.0011 0.002
?=0, 0.0011 Emp ?, 0.0007
Approx 1M
Empirical 10k
App ?
?
?=1
Emp ?

?=0.9
?=0.8
?=0.7
?=0.6
Return Length, n Empirical 1M
?=0.5
-0.5
?=0.4
Approx 10k
?=0.3
Approx 1M
Empirical 10k
?=0.2
10
Return Length Empirical 1M
?=0
0
IS
0
?=0.1
0 WIS
Variance
0.16
Return Type
Approx 10k
(e) Approximate and empiri- (f) Approximate and empirical weights (g)
Mean squared error from 10000 tracal diagonals of ?. for each return. jectories.

Figure 3: Digital Marketing Results.

40
30
40
25
30
20
30
20
20
10
20
10
0
10
0 0
?10 0
0 0
15 10

5
5
20
10
25
25
15
15
10 20
20
10
15
15
5
5
20
10
10 20
5
5
15
15
10 20
20
10
15
15
5 0 0
10 20
5
5 25
25

(a) Empirical ? from (b) Empirical ? from 10 (c) Approximate ? from (d) Approximate ? from 10000 trajectories. trajectories. 10000 trajectories. 10 trajectories. 1

Approx 10K
Empirical 10
Approx 10
App ?
?
Emp ?
?=0.6
?=0.5
?=0.4
?=0.3
Return Length, n Empirical 10K

?=0

-0.5

?=0.2

Approx 10

?=0.1

Approx 10K

Empirical 10

IS

Return Length Empirical 10K

20

WIS

10

?=1

?=1, 3.47436 App ?, 3.1070

1 0

?=0.9

20

0

?=0.8

1

?=0, 3.2102

10

?=0.7

20

Mean Squared Error

Weight, w(n, 10)

Variance

40

100

Return Type

(e) Approximate and empiri- (f) Approximate and empirical weights (g)
Mean squared error from 10 trajeccal diagonals of ?. for each return. tories.

Figure 4: Functional Electrical Stimulation Results.

6

(i)

(L)

and Cov[Rst , Rst ]. 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. Estimating Diagonal Entries of ?

4.2

The remaining question is how best to approximate the diagonal of ? 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 ? 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. (i) Let vi denote the sample variance of Rst for i = 1 . . . L. Let v+ be the largest sample variance: v+ = maxi?{1,...,L} vi . We parameterize the diagonal using four parameters, k1 , k2 , v+ , and vL : ? ? if i = 1 ?k 1 ? if i = L ?k1 ,k2 ,v+ ,vL (i, i) = vL ? ?min{v , k k (1?t) } otherwise. + 1 2 ?(1, 1) = k1 sets the initial variance, and vL is the variance of the Monte Carlo return. The parameter v+ enforces a ceiling on the variance of the i-step return, and k2 captures the growth rate of the variance, much like ?. We select the k1 and k2 that minimize the mean squared error between ? i) and vi , and set v+ and vL directly from the data.2 ?(i, This reduces the problem of estimating ?, an L ? L matrix, to estimating four numbers from return ? as computed from many trajectories. data. Consider Figures 1(c), 2(c), 3(c), and 4(c), which depict ? The differences between these estimates and the ground truth show that this parameterization is not perfect, as we cannot represent the true ? 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 ? from very few trajectories. Figures ? when computed from only a few trajectories. Note their similarity 1(d), 2(d), 3(d), and 4(d) show ? ? to ? 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 ?-return is provided in Algorithm 1. Unlike the ?-return, which can be computed from a single trajectory, the ?-return requires a set of trajectories in order to estimate ?. The pseudocode assumes that every trajectory is of length L, which can be achieved by padding shorter trajectories with absorbing states.

2

We include the constraints that k2 ? [0, 1] and 0 ? k1 ? v+ .

7

Algorithm 1: Computing the ?-return. Require: n trajectories beginning at s and of length L. (i)

1. Compute Rs for i = 1, . . . , L and for each trajectory. (i)

2. Compute the sample variances, vi = Var(Rs ), for i = 1, . . . , L. 3. Set v+ = maxi?{1,...,L} vi . 4. Search for the k1 and k2 that minimize the mean squared error between vi and ? k ,k ,v ,v (i, i) for i = 1, . . . , L. ? 1

2

+

L

? k ,k ,v ,v (i, i), using the 5. Fill the diagonal of the L ? L matrix, ?, with ?(i, i) = ? 1 2 + L optimized k1 and k2 . 6. Fill all of the other entries with ?(i, j) = ?(i, i) where j ¿ i. If (i = L or j = L) and i 6= j then set ?(i, j) = 0 instead. 7. Compute the weights for the returns according to (3). 8. Compute the ?-return for each trajectory according to (1).

15

# 5

## Experiments

Approximations of the ?-return could, in principle, replace the ?-return in the whole family of TD(?) algorithms. However, using the ?-return for TD(?) raises several interesting questions that are beyond the scope of this initial work (e.g., is there a linear-time way to estimate the ?-return? Since a different ? is needed for every state, how can the ?-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.3 Notice that, for all domains, using the ?-return (the EMP ? and APP ? labels) results in lower MSE than the ?-return and the ?-return with any setting of ?.

# 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 ? and ?-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 ?-return and TD(?) 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 ?-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 ?-return that allows for value function approximation with the same time complexity as TD(?), or which better leverages our knowledge that the environment is Markovian? Would TD(?) using the ?-return be convergent in the same settings as TD(?)? 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. 3 To compute the MSE we used a large number of Monte Carlo rollouts to estimate the true value of each policy.

8

# 2　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? : 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. Theocharous 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. Theocharous, 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 TwentyFirst Innovative Applications of Artificial Intelligence, pages 165?172, 2009. [15] P. M. Pilarski, 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.

9