

Random-Sampling Monte-Carlo Tree Search Methods for Cost Approximation in Long-Horizon Optimal Control

Shankarachary Ragi, *IEEE Senior Member* and Hans D. Mittelmann

Abstract—In this paper, we develop Monte-Carlo based heuristic approaches to approximate the objective function in long horizon optimal control problems. In these approaches, to approximate the expectation operator in the objective function, we evolve the system state over multiple trajectories into the future while sampling the noise disturbances at each time-step, and find the average (or weighted average) of the costs along all the trajectories. We call these methods *random sampling - multipath hypothesis propagation* or RS-MHP. These methods (or variants) exist in the literature; however, the literature lacks results on how well these approximation strategies converge. This paper fills this knowledge gap to a certain extent. We derive convergence results for the cost approximation error from the RS-MHP methods and discuss their convergence (in probability) as the sample size increases. We consider two case studies to demonstrate the effectiveness of our methods - a) linear quadratic control problem; b) UAV path optimization problem.

Index Terms—Long horizon optimal control, cost approximation, approximate dynamic programming, multipath hypothesis propagation.

I. INTRODUCTION

Long-horizon optimal control problems appear naturally in robotics, advanced manufacturing, and economics, especially in applications requiring decision making in stochastic environments. Often these problems are solved via dynamic programming (DP) formulation [2]. DP problems are notorious for their computational complexity, and require approximation approaches to make them tractable. A plethora of approximation techniques called *approximate dynamic programs* (ADPs) exist in the literature to solve these problems approximately. Some of the commonly used ADPs include *policy rollout* [3], *hindsight optimization* [4], [5], etc. A survey of the ADP approaches can be found in [2]. Feature-based techniques and deep learning methods are gaining importance in the development of ADP approaches as discussed in [6]. These approximation techniques have

This work was supported in part by Air Force Office of Scientific Research under grant FA9550-19-1-0070. This paper was presented in part at The 4th IEEE Conference on Control Technology and Applications 2020 [1].

Shankarachary Ragi (corresponding author) is with Department of Electrical Engineering, South Dakota School of Mines and Technology, Rapid City, SD 57701, USA shankarachary.ragi@sdsmt.edu

Hans D. Mittelmann is with the School of Mathematical and Statistical Sciences, Arizona State University, Tempe, AZ 85281, USA nlanchie@asu.edu, mittelmann@asu.edu

been successfully adopted to solve real-time problems such as a UAV guidance control problem in [7]–[9]. Certain ADP approaches, especially the methods based on approximation in value space, require numerical approximation of the expectation in the objective function [7]. In this study, our objective is to develop Monte-Carlo-based approaches to approximate the expectation in the objective function in the long (but finite) horizon optimal control problems, and study their convergence. A preliminary version of the parts of this paper were published as [1]. This paper differs from the conference paper [1] in the following ways: 1) we include detailed proofs omitted in the conference version; 2) we derive new convergence results and proofs in Section II-A; 3) we implement our methods for a new case study - UAV path optimization problem.

A. Preliminaries

A long horizon optimal control problem is described as follows. Let x_k be the state vector for a system at time k , which evolves according to a discrete stochastic process as follows:

$$x_{k+1} = f(x_k, u_k, w_k) \quad (1)$$

where $f(\cdot)$ represents the state-transition mapping, u_k is the control vector, and w_k random disturbance. Let $g(x_k, u_k)$ represent the cost (a real value) of being in state x_k and performing action u_k . The functions f and g are independent of k in our study, but can generally depend on k . The goal is to optimize the control vectors $u_k, k = 0, \dots, H-1$ such that the expected cumulative cost is minimized, i.e., the goal leads to solving the following optimization problem

$$\min_{u_k, k=0, \dots, H-1} \mathbb{E} \left[\sum_{k=0}^{H-1} g(x_k, u_k) \right], \quad (2)$$

where H is the length of the planning horizon. Let x_0 be the initial state and according to the dynamic programming formulation the optimal cost function is given by

$$J_0^*(x_0) = \min_{u_0} \mathbb{E} [g(x_0, u_0) + J_1^*(x_1)], \quad (3)$$

where J_1^* represents the optimal cost-to-go from time $k = 1$, and $x_1 = f(x_0, u_0, w_0)$. In this study, *long horizon* refers to the condition that H is sufficiently large that the

optimal policy is approximately *stationary* (independent of k). Solving the above optimization problem is not tractable mainly due to two reasons: the expectation $E[\cdot]$ and the optimal cost-to-go J_1^* are hard to evaluate and are usually approximated by numerical methods or ADP approaches.

An ADP approach called *nominal belief-state optimization* (NBO) [7], [10] was developed primarily to approximate the above expectation. In NBO, the expectation is replaced by a sample state trajectory generated with an assumption that the future noise variables in the system take so called nominal or mean values, thus making the above objective function deterministic. The NBO method was developed to solve a UAV path optimization problem, which was posed as a *partially observable Markov decision process* (POMDP). POMDP generalizes the long horizon optimal control problem described in Eq. 2 in that the system state is assumed to be “partially” observable, which is inferred via using noisy observations and Bayes rules. Although the performance of the NBO approach was satisfactory, in that it allowed to obtain reasonably optimal control commands for the UAVs, it ignored the uncertainty due to noise disturbances thus leading to inaccurate evaluation of the objective function. To address this challenge, certain methods exist in the literature usually referred to as Monte-Carlo Tree Search (MCTS) methods as surveyed in [11].

Inspired from the NBO method and MCTS methods, we develop a new MCTS method called *random sampling - multipath hypothesis propagation* (RS-MHP) and derive convergence results. In this study, we use the NBO approach as a benchmark for performance assessment since RS-MHP builds on the NBO approach.

II. RANDOM SAMPLING MULTIPATH HYPOTHESIS PROPAGATION (RS-MHP)

In the NBO method, the expectation is replaced by a sample trajectory of the states (as opposed to random states) generated by

$$\tilde{x}_{k+1} = f(\tilde{x}_k, u_k, \bar{w}_k), k = 0, \dots \quad (4)$$

where $\tilde{x}_0 = x_0$ (initial state or current state), and \bar{w}_k is the mean of the random variable w_k . Thus, the long horizon optimal control problem, with NBO approximation, reduces to

$$\min_{u_k} \sum_{k=0}^{H-1} g(\tilde{x}_k, u_k). \quad (5)$$

The above reduced problem, without the need for evaluating the expectation, can significantly reduce the computational burden in solving the long horizon control problems. However, the downside with this approach is it completely ignores the uncertainty in the state evolution, and may generate severely sub-optimal controls. To

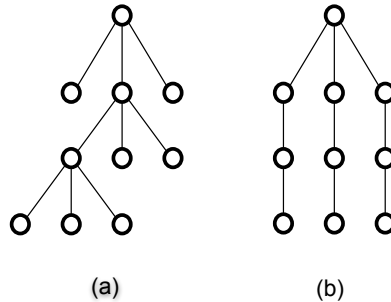


Figure 1. State trajectory sampling models: (a) tree branching model, (b) non-overlapping branching model.

overcome this trivialization, we develop a Monte-Carlo approach to approximate the expectation described as follows. We will follow the tree-like sampling approach as in Figure 1(a). For time step $k = 1$, we sample the probability distribution of the noise disturbance N times to generate the samples w_0^i with corresponding probability p_0^i , $i = 1, \dots, N$. Using these, we generate N sample states at $k = 1$ generated according to

$$x_1^i = f(x_0, u_0, w_0^i), \forall i. \quad (6)$$

We repeat this sampling approach for time $k = 2$, i.e., we generate N noise samples w_1^i with corresponding probability p_1^i , $i = 1, \dots, N$. Using these noise samples and the sample states from the previous time step, we generate N^2 sample states at $k = 2$ according to

$$x_2^{i,j} = f(x_1^i, u_1, w_1^j), \forall i, j. \quad (7)$$

We repeat the above sampling procedure until the last time step $k = H - 1$ to generate N^{H-1} possible state evolution trajectories using N noise samples generated in each time step as depicted in Figure 1(a). Sampling approach in Figure 1(b) will be discussed later.

One can now replace the expectation in Eq. 2 with the weighted average of the cumulative cost corresponding to each state evolution trajectory, where the weights are the probabilities or likeliness of the trajectories. Clearly, the number of possible state trajectories grow exponentially with the horizon length H . Although this approach is not novel as many such methods exist in the literature often classified as Monte-Carlo Tree Search methods, our study is focused on deriving convergence results of RS-MHP approaches.

To avoid the exponential growth in our RS-MHP approach, at each time step we retain only M sample states and prune the remaining states, and if the number of sample states at a given time instance is less than or equal to M , we do not perform pruning. For pruning, at each time k , we rank the state trajectories up to time k according to their likeliness (obtained by multiplying the

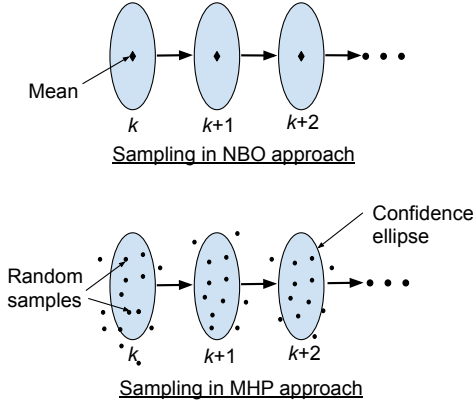


Figure 2. Sampling probability distributions of noise variables: NBO vs. MHP.

probabilities of all the noise samples that generated the trajectory) and retain the top M trajectories with highest likelihood and prune the rest. With this procedure, at $k = H - 1$, there would be only M state trajectories. With pruning, the number of trajectories remains a constant irrespective of the time horizon length. An illustration of the above RS-MHP approach is shown in Figure 2 along with the NBO approach. Here, we consider pruning based on likelihood of the state trajectories as the costs from these trajectories have higher contribution in the cost function in Eq. 1 than the less likely trajectories. We will consider other pruning strategies to further improve the approximation error in our future study.

Let $i = 1, \dots, M$ represent the indices of the M distinct state trajectories with q_1, q_2, \dots being their likelihood index evaluated using the probabilities of the noise samples that generate the trajectory i over time. Let J represent the actual objective function as described below

$$J = \mathbb{E} \left[\sum_{k=0}^{H-1} g(x_k, u_k) \right]. \quad (8)$$

We can now approximate the objective function J in four possible ways as described below (assuming $N > M$). Let x_k^i represent the state at time k in the i th state trajectory.

(I) *Sample Averaging.* We can simply approximate the expectation with an average over all possible trajectories as follows:

$$\begin{aligned} \text{No pruning: } J &\approx \bar{J}_{NP} = \frac{1}{N^{H-1}} \sum_{i=1}^{N^{H-1}} \left(\sum_{k=0}^{H-1} g(x_k^i, u_k) \right) \\ \text{With pruning: } J &\approx \bar{J}_P = \frac{1}{M} \sum_{i=1}^M \left(\sum_{k=0}^{H-1} g(x_k^i, u_k) \right) \end{aligned} \quad (9)$$

(II) *Weighted Sample Averaging.* We can also approximate the expectation with a weighted average with

weights being the normalized likelihood indices of the state trajectories given by $q_i, i = 1, \dots$ (and \bar{q}_i in the pruned case) as follows:

$$\begin{aligned} \text{No pruning: } J &\approx \bar{J}_{NP} = \frac{1}{N^{H-1}} \sum_{i=1}^{N^{H-1}} q_i \left(\sum_{k=0}^{H-1} g(x_k^i, u_k) \right) \\ \text{With pruning: } J &\approx \bar{J}_P = \frac{1}{M} \sum_{i=1}^M \bar{q}_i \left(\sum_{k=0}^{H-1} g(x_k^i, u_k) \right). \end{aligned} \quad (10)$$

where $\sum_{i=1}^{N^{H-1}} q_i = N^{H-1}$ and $\sum_{i=1}^M \bar{q}_i = M$.

For a given sequence of control decisions u_0, u_1, \dots , let g_i denote the cost of the i th trajectory given by

$$g_i = \sum_{k=0}^{H-1} g(x_k^i, u_k). \quad (11)$$

Clearly, g_1, g_2, \dots are identically distributed random variables, but are dependent due to the overlapping state trajectories in the tree-like sampling approach in Figure 1(a), where $\mathbb{E}[g_i] = J, \forall i$.

The below result suggests that with sufficient number of sample state trajectories (large N), the approximation error in \bar{J}_{NP} becomes small enough to ignore.

Proposition 2.1: For any given sequence of actions u_0, u_1, \dots , if the random variables g_1, g_2, \dots have finite variances, \bar{J}_{NP} converges to J in probability.

Proof: From Ex. 254 in [12], we know that $\bar{J}_{NP} \xrightarrow{P} J$ if

$$\lim_{|i-j| \rightarrow \infty} \text{Cov}(g_i, g_j) = 0, \quad (12)$$

where $\text{Cov}()$ represents covariance. Suppose, the sequence g_1, g_2, \dots is arranged such that g_1 represents the cost for the left-most branch in Figure 1(a), and g_2 representing the second branch from the left, and so on. Clearly, the first g_1, g_2, \dots, g_N are dependent random variables as they share the same parent node, whereas the next N terms $g_{N+1}, g_{N+2}, \dots, g_{2N}$, although dependent among themselves, are independent of the previous N terms (as these branches evolve from a separate parent node), and so on. Thus, $\text{Cov}(g_i, g_j) = 0$ if $|i - j| > N$, which implies $\lim_{|i-j| \rightarrow \infty} \text{Cov}(g_i, g_j) = 0$. ■

Furthermore, we can apply similar arguments to prove the convergence of \bar{J}_{NP} in probability.

Proposition 2.2: For a given sequence of actions u_0, \dots, u_{H-1} , if g_1, g_2, \dots have finite variances, then \bar{J}_{NP} converges to J in probability.

Proof: From [13], we know that if $\bar{J}_{NP} \xrightarrow{P} J$ (which is true as shown in Proposition 2.1), and if the weights q_1, q_2, \dots are monotonically decreasing, then $\bar{J}_{NP} \xrightarrow{P} J$. Without loss of generality, we can arrange the trajectory costs g_i such that their likelihood indices are monotonically decreasing, i.e., $q_1 \geq q_2 \geq q_3 \geq \dots$, which completes the proof. ■

A. Non-overlapping State Trajectories or Tree Branches

Suppose the state sample trajectories are generated independently of each other, where the state trajectories do not share any common state samples as depicted in Figure 1b. In this new sampling approach, given u_0, u_1, \dots are the control decisions over the planning horizon, let p_i represent the cost associated with the i th state trajectory. We can approximate the LHC objective function as follows:

$$\begin{aligned}\bar{J}_N &= \frac{1}{N} \sum_{i=1}^N p_i \\ \tilde{J}_N &= \frac{1}{N} \sum_{i=1}^N q_i p_i,\end{aligned}\quad (13)$$

where q_i represents the likeliness index of the i th trajectory and $\sum_i q_i = N$. From propositions 2.1 and 2.2, we can verify that $\bar{J}_N \xrightarrow{P} J$ and $\tilde{J}_N \xrightarrow{P} J$. Furthermore, since p_1, p_2, \dots are i.i.d., due to the strong law of large numbers, we can verify that \bar{J}_N converges to J almost surely. We can further derive the rate of convergence (in probability) for a special case as discussed below. Suppose the state-transition and cost functions are linear (motivated by the fact that the linear models capture the state dynamics well in most control problems) as described below:

$$\begin{aligned}x_{k+1} &= Ax_k + Bu_k + w_k, w_k \sim \mathcal{N}(0, \Sigma) \\ g(x_k, u_k) &= Cx_k + Du_k,\end{aligned}\quad (14)$$

where $g(x_k, u_k)$ is a scalar function. The cost from the sample trajectory i is given by

$$p_i = \sum_{k=1}^H g(x_k^i, u_k) = \sum_{k=1}^H (Cx_k^i + Du_k), \quad (15)$$

where x_k^i is the sampled state at time step k from the i th trajectory. Using the linear expressions in Eq. 14, we can verify p_i further satisfies the following equation:

$$p_i - \mathbb{E}[p_i] = C \left[\sum_{k=0}^{H-1} \left(\sum_{q=0}^{H-k-1} A^q \right) w_k \right] = C \left[\sum_{k=0}^{H-1} \mathcal{A}_k w_k \right], \quad (16)$$

where $\mathcal{A}_k = \sum_{q=0}^{H-k-1} A^q$.

Proposition 2.3: For a given sequence of actions u_0, \dots, u_{H-1}

$$\mathbb{P}(|J_N - J| \geq \varepsilon) \leq \frac{\text{constant}}{N\varepsilon^2}. \quad (17)$$

Proof: Let p represent the cost for a sampled state trajectory. Using Eq. 16, we can verify

$$\begin{aligned}\text{Var}(p) &= \mathbb{E}[(p - \mathbb{E}[p])^T (p - \mathbb{E}[p])] \\ &= C \left[\sum_{k=0}^{H-1} \mathcal{A}_k \Sigma \mathcal{A}_k^T \right] C^T,\end{aligned}\quad (18)$$

which is a real scalar. Thus, $\text{Var}(J_N) = \text{Var}(p)/N$.

Using Chebyshev's inequality, we can verify easily that

$$\mathbb{P}(|J_N - J| \geq \varepsilon) \leq \frac{\text{Var}(p)}{N\varepsilon^2} = \frac{C \left[\sum_{k=0}^{H-1} \mathcal{A}_k \Sigma \mathcal{A}_k^T \right] C^T}{N\varepsilon^2}. \quad (19)$$

Furthermore,

$$\lim_{N \rightarrow \infty} \mathbb{P}(|J_N - J| \geq \varepsilon) = 0, \quad (20)$$

which shows the convergence in probability as well.

III. CASE STUDIES

We implement the above-discussed MHP methods in the context of two case studies: (a) linear quadratic Gaussian control (LQG); (b) path planning for unmanned aerial vehicles (UAVs). These case studies are discussed below.

A. Linear Quadratic Problem

Although there are closed-form solutions for LQG problems, the below example allows us to quantify the benefits of using RS-MHP methods over existing similar methods, particularly NBO. Let the system state evolve according to the following linear equation:

$$x_{k+1} = (1-a)x_k + au_k + w_k, \quad w_k \sim \mathcal{N}(0, \sigma^2), \quad (21)$$

where $0 < a < 1$ is a constant, and w_k is a random disturbance modeled by a zero-mean Gaussian distribution with variance σ^2 . The cost function over the time-horizon H is defined as follows:

$$J = \mathbb{E} \left[r(x_H - T)^2 + \sum_{k=0}^{H-1} u_k^2 \right], \quad (22)$$

where r and T are constants. This is a simplified oven temperature control example borrowed from [14].

If we apply the traditional NBO method, assuming $H = 2$, the cost function J is approximated (assuming nominal values or zeros for w_0 and w_1) as

$$J_{\text{NBO}} = r((1-a)^2 x_0 + a(1-a)u_0 + au_1 - T)^2 + u_0^2 + u_1^2 \quad (23)$$

and the exact cost function J can be evaluated analytically as

$$\begin{aligned}J &= r((1-a)^2 x_0 + a(1-a)u_0 + au_1 - T)^2 + u_0^2 + u_1^2 \\ &\quad + r\sigma^2((1-a)^2 + 1).\end{aligned}\quad (24)$$

We notice the approximation error due to the NBO method is

$$|J_{\text{NBO}} - J| = r\sigma^2((1-a)^2 + 1). \quad (25)$$

This approximation error for a generic time-horizon H (the above error term is derived for $H = 2$) is given by

$$|J_{\text{NBO}} - J| = r\sigma^2 \sum_{n=0}^{H-1} (1-a)^{2n}. \quad (26)$$

The above expression suggests that the NBO approximation error can be significantly high depending on the parameters a , σ , and r . With MHP approximation, the cost function reduces to

$$J_{\text{MHP}} = \frac{1}{P} \left(\sum_{i=1}^P r(x_H^i - T)^2 \right) + \sum_{k=0}^{H-1} u_k^2, \quad (27)$$

where P is the number of state-trajectories generated using the MHP approach, and x_H^i is the final state in the i th trajectory. Lemma 2.1 shows that the approximation error due to the above MHP method converges (in probability) to zero. We verify this result with a numerical simulation, where we implement the NBO and the MHP methods with the following assumptions: $x_0 = 0, r = 10, T = 1, H = 2, u_0 = 0.55, u_1 = 0.17, \sigma = 1$. We vary P from 100 to 10000 with increments of 100. Figure 3 shows the cost function approximated using MHP and NBO methods. The figure clearly demonstrates that the error due to NBO approximation can be significantly high, while MHP performs better in cost approximation.

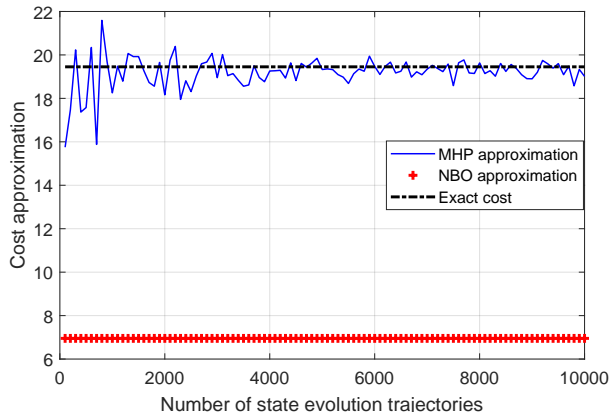


Figure 3. LQG problem: MHP vs. NBO

B. UAV path planning problem

We consider a UAV path planning problem, where the goal is to optimize the kinematic controls of a UAV to maximize a target tracking performance measure. Here, the UAV is assumed to be equipped with a sensor on-board that generates the location measurements of the target (a ground-based moving vehicle) corrupted by random noise. A detailed description of the problem can be found in [7]. In [7], we posed this problem as a *partially observable Markov decision process* (POMDP), where the POMDP led to solving a long horizon optimal control problem. We applied the NBO approach to solve

the above POMDP. The resulting UAV path optimization problem is summarized as follows:

$$\min_u \mathbb{E} \left[\sum_{k=0}^{H-1} \text{tr}(\mathbf{P}_k(u)) \right] \xrightarrow{\text{NBO approx.}} \min_u \sum_{k=0}^{H-1} \text{tr}(\hat{\mathbf{P}}_k(u)),$$

where $\mathbf{P}_k(u)$ (a random variable) represents the error covariance matrix corresponding to the state of the system, $\text{tr}()$ represents the matrix trace operator, u is the sequence of UAV kinematic controls (e.g., forward acceleration and bank angle) applied over the discrete time planning horizon of length H steps. After NBO approximation, the expectation over the random evolution of $\mathbf{P}_k(u)$ is replaced with the nominal sequence of the state covariance matrices $\text{tr}(\hat{\mathbf{P}}_k(u))$.

We now approximate the above objective function using the RS-MHP approach as follows:

$$\min_u \mathbb{E} \left[\sum_{k=0}^{H-1} \text{tr}(\mathbf{P}_k(u)) \right] \xrightarrow{\text{RS-MHP approx.}} \min_u \frac{1}{N_T} \sum_{i=1}^N \sum_{k=0}^{H-1} \text{tr}(\tilde{\mathbf{P}}_k^i(u)),$$

where $\tilde{\mathbf{P}}_k^i$ represents the state covariance matrix obtained from the i th state trajectory generated from the RS-MHP approach, and N_T is the number of state trajectories. We implement this RS-MHP approach in MATLAB and run a Monte-Carlo study to see the impact of N_T on the performance of the above UAV path planning algorithm, which is measured by the average target location estimation error. Figure 4 shows the cumulative distribution of average target location estimation errors from the RS-MHP approach with $H = 6$, and for N_T set to 50, 100, and 250. The figure shows a gradual increase in the UAV path optimization performance with increasing N_T as expected. This result, as expected, also suggests that pruning methods (discussed in the previous section) would degrade the performance of the RS-MHP methods but can provide gains in terms of computational intensity.

RS-MHP has better capability in approximating the expectation operator in Eq. 1 than the NBO approach as we consider multiple hypotheses of state trajectories in RS-MHP as opposed to a single hypothesis in NBO as demonstrated in Figure 5. This is demonstrated in the above case studies.

IV. CONCLUSIONS

In this paper, we developed a Monte-Carlo tree search method called *random sampling - multipath hypothesis propagation* or RS-MHP to approximate the expectation operator in long horizon optimal control problems. Although variants of these methods exist in the literature, we focused on the convergence analysis of these approximation methods. The basic theme of these methods

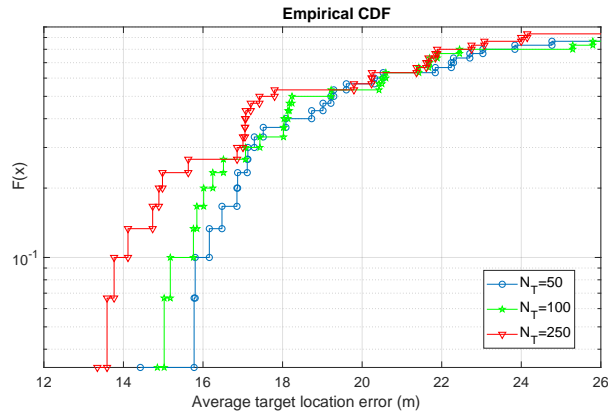


Figure 4. Cumulative distribution of average target location errors. Here N_T represents the number of state evolution trajectories.

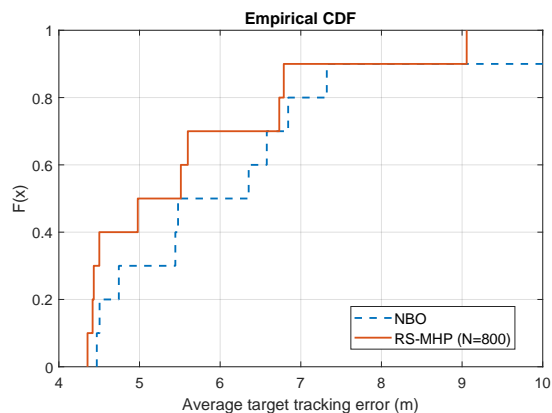


Figure 5. Cumulative distribution of average target location errors: NBO vs. RS-MHP.

is to evolve the system state over multiple trajectories into the future while sampling the noise disturbances at each time-step. We derive convergence results that show that the cost approximation errors from our RS-MHP methods converge (in probability) toward zero as the sample size increases. We conducted a numerical study to assess the performance of our methods in two case studies: linear quadratic control problem and UAV path optimization problem. In both case studies, we demonstrated the benefits of our approach against an existing approach called *nominal belief-state optimization* or NBO (used as a benchmark).

V. ACKNOWLEDGMENT

The authors would like to thank Nicolas Lanchier, Arizona State University, for his valuable inputs and feedback on the convergence results discussed in this paper.

REFERENCES

- [1] S. Ragi and H. D. Mittelmann, "Random-sampling multipath hypothesis propagation for cost approximation in long-horizon optimal control," in *Proc. 2020 IEEE Conference on Control Technology and Applications (CCTA)*, Montreal, Canada, 2020, pp. 14–18.
- [2] E. K. P. Chong, C. M. Kreucher, and A. O. Hero, "Partially observable Markov decision process approximations for adaptive sensing," *Discrete Event Dynamic Systems*, vol. 19, no. 3, pp. 377–422, Sep 2009.
- [3] D. P. Bertsekas and D. A. Castanon, "Rollout algorithms for stochastic scheduling problems," *J. Heuristics*, vol. 5, pp. 89–108, 1999.
- [4] E. K. P. Chong, R. L. Givan, and H. S. Chang, "A framework for simulation-based network control via hindsight optimization," in *Proc. 39th IEEE Conf. Decision and Control*, Sydney, Australia, 2000, pp. 1433–1438.
- [5] G. Wu, E. K. P. Chong, and R. Givan, "Burst-level congestion control using hindsight optimization," *IEEE Trans. Autom. Control*, vol. 47, pp. 979–991, 2002.
- [6] D. Bertsekas, "Feature-based aggregation and deep reinforcement learning: A survey and some new implementations," *IEEE/CAA Journal of Automatica Sinica*, no. 1, 2019.
- [7] S. Ragi and E. K. P. Chong, "UAV path planning in a dynamic environment via partially observable Markov decision process," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 49, pp. 2397–2412, 2013.
- [8] —, "Dynamic UAV path planning for multitarget tracking," in *Proc. American Control Conf.*, Montreal, Canada, 2012, pp. 3845–3850.
- [9] S. Ragi and H. D. Mittelmann, "Mixed-integer nonlinear programming formulation of a UAV path optimization problem," in *Proc. American Control Conf.*, Seattle, WA, 2017, pp. 406–411.
- [10] S. Miller, Z. Harris, and E. K. P. Chong, "A POMDP framework for coordinated guidance of autonomous UAVs for multitarget tracking," *EURASIP Journal on Advances in Signal Processing*, 2009.
- [11] C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A survey of Monte Carlo tree search methods," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 4, no. 1, pp. 1–43, March 2012.
- [12] T. Cacoullos, *Exercises in Probability*. New York: Springer-Verlag, 1989.
- [13] N. Etemadi, "Convergence of weighted averages of random variables revisited," *Proc. American Mathematical Society*, vol. 134, no. 9, pp. 2739–2744, 2006.
- [14] D. P. Bertsekas. Lecture on reinforcement learning and optimal control. [Online]. Available: http://www.mit.edu/~dimitrib/Slides_Lecture2_RLOC.pdf