

A Novel Approach to Reproduced Kernel Hilbert Space for Artificial Intelligence Control: Using Monte-Carlo Estimates of Operators Arising

Shima Asghari [†]*, Nahid Khanlari ^{††}

Department of Applied Mathematics, Islamic Azad University, Hamedan Branch, Hamedan 65138, Iran

Abstract—An embedding of stochastic optimal control problems of artificial intelligence form into reproducing kernel Hilbert spaces is presented in this study. A model-free, non-parametric approach for calculation of an approximate solution to the control problem is obtained by consistent, sample based estimates of the embedding. The solve the problem; it decomposes into two components; invariant and task dependent. Hence, the solution presented in the current paper used sample data more efficiently than previous sample based approaches in this field by some innovations such as allowing sample re-use across tasks. To show the efficiency of the introduced approach, numerical examples on test problems are presented.

Index Terms— Monte-Carlo Estimates, Operators, Artificial Intelligence, Hilbert Space, Reproducing Kernel, Double Slit, Applied Mathematics

1 Introduction

In spite of challenges in solving general non-linear stochastic optimal control (SOC) and Reinforcement Learning (RL) problems, a type of problems which can be solved by closed form solutions are recently identified [1-8]. To solve these problems, it is necessary to evaluate an artificial intelligence, which is equivalent to evaluating a partition function which in turn is a tough problem. However, these solutions are interesting due to their several practical applications, resulting from the possibility of

applying Monte-Carlo and variational methods for solving these problems [9, 10]. Analytical evaluating the required artificial intelligence, based on linear operators acting on state vectors, is possible in some special cases such as linear dynamics and quadratic costs [11, 14, 16, and 17]. It is shown in the current paper that the artificial intelligence can be evaluated, regarding the covariance operators acting on elements of the Hilbert space, by appropriate embedding of it into a reproducing kernel Hilbert space (RKHS). Although a tractable solution to the SOC problem does not yield in itself, consistent estimators of the required operators lead to efficient non-parametric algorithms [11-23].

A critical prominence of estimating the operators other than directly estimating the artificial intelligence, as the goal of the previous applications of Monte-Carlo methods, is that some deficiencies of previous methods can be eliminated without considerable loss of their advantages [24]. In this approach, the complexity of sample is considerably reduced due to separating the problem into an invariant and task varying

- *Shima Asghari, Corresponding Author, Ph.D. Student in Applied Mathematics, Department of Applied Mathematics, Islamic Azad University, Hamedan Branch, Hamedan 65138, Iran.*
- *Nahid Khanlari, Ph.D. Student in Applied Mathematics, Department of Applied Mathematics, Islamic Azad University, Hamedan Branch, Hamedan 65138, Iran.*

component [25]. It allows efficient sample re-use across tasks and leads to a form of transfer learning which in turn, leads to the situation where any change in the task including, for e.g., different start states, necessitate acquiring new samples [26-39]. Moreover, as the approach remains model-free, it is applicable to the RL setting, unlike to variational [40-45] or function approximation [46-59] approaches. As a result, it is further distinguished through convergence guarantees. The operators becomes state-dimensionality independent by embedding of RKHS. It leads to better scalability of operators. However, by informing choices of sampling procedures and kernel, incorporating the prior knowledge about both tasks and dynamics is effectively possible [60].

It should be noted that the presented approach in this study is not limited to problems which are in context of artificial intelligence stochastic optimal control. For solving linearly solvable MDPs by [61, 62], inference control by [63-73] and free energy control of [74, 75], an underlying problem of equivalent form should be solved. As a result, the proposed methods are directly applicable in these fields. In addition, a generalization of artificial intelligence control to develop an optimal policy for general SOC problems is described by [76-83]. Moreover, [84] and [85] were proposed artificial intelligence formulations for discounted and average cost infinite horizon problems and risk sensitive controls, respectively, while in this study, finite horizon problems are critically focused.

A brief review of the formulation of a SOC problem in the artificial intelligence control framework of [86-90], which has two specific properties related to the presented approach in this study. The first of these characteristics is that this formulation solves the optimal policy according to a state desirability function Ψ in a closed form manner. The second one is the possibility of expressing this function Ψ as a conditional expectation of the product of an immediate cost and a future desirability. Therefore, the desirability function can be computed by identifying this expression with an inner product in a RKHS. The

concept is formulated, as the model-based case is considered and the analytical form of the RKHS based evaluation also is extracted. Discussing about the possible estimating of this operator from transition samples, which leads to performing in form of a finite dimensional inner product, considering the model-free case. For solutions of the SOC problem are summarized as (i) transition samples gathered and then, embedding operator estimated, (ii) the desirability function and immediate cost represented as elements of RKHSs, (iii) the desirability function evaluated, recursively, as a series of inner products based on the estimated operator –including a series of matrix multiplications – and (iv) the optimal controls computed based on the obtained desirability function. The methodology is developed, accompanying by a series of alternative estimators so that the goal is either reducing the computational costs or allowing more efficient use of the sample data. The verifications of the proposed methods by experimental data.

2 General Description of Artificial Intelligence Control

The artificial intelligence approach to stochastic optimal control, proposed by [79-84] (see also [85-90]), is reviewed in the current section.

Considering $x \in R^{d_x}$ as the system state and $u \in R^{d_u}$ as the control signals, a continuous time stochastic system can be taken into account as the following form:

$$dx = f(x, t)dt + B(x, t)(udt + d\zeta) \quad (1)$$

where $d\zeta$ is a multivariate Wiener process with $E[d\zeta^2] = Q(x, t)dt$, and f , B and Q may be non-linear functions. It is worthwhile to note that the system is linear in the controls while both noise and controls act in the same subspace. By an objective of the following form, the best Markov policy, i.e., $u(t) = \pi(x(t), t)$ is searched:

$$J^\pi(x, t) = E_{X^\pi(0)|x} \left[C \cdot (X^\pi(T)) + \int_t^T C(X^\pi(s), s) + u(s)^T H u(s) ds \right] \quad (2)$$

where T is a specified terminal time, $C \cdot$ and C are a terminal and running cost, respectively. The expectation is taken w.r.t. $X^\pi(0)$, the paths of Eq. (1) starting in x and following policy π . By requiring the quadratic control cost, as $H \in R^{d_u \times d_u}$, to satisfy $Q = \lambda B H^{-1} B^T$ for some constant scalar $\lambda > 0$, it is further constrained.

It was shown by [79-90] that the optimized objective for this form of problem can be stated as:

$$J^*(x, t) = \min_{\pi} J^\pi(x, t) = -\lambda \log \Psi(x, t) \quad (3)$$

where Ψ is a state desirability function and can be obtained by the following path integral:

$$\Psi(x, t) = E_{X^0(0)|x} \left[e^{-\int_t^T \frac{1}{\lambda} C(X^0(s), s) ds} \Psi(X^0(T), T) \right] \quad (4)$$

with $\Psi(0, T) = \exp\{-C \cdot(0) / \lambda\}$. The expectation in Eq. (4) is taken w.r.t. uncontrolled path of the dynamics Eq. (1), i.e. those under the policy $\pi^0(0, 0) = 0$, starting in x_t .

It can be described the optimal policy $\pi^*(x, t)$, directly, in terms of as the following due to linear control with quadratic control cost and Ψ :

$$\pi^*(x, t) = -H^{-1} B(x)^T \nabla_x J^*(x, t) \quad (5)$$

$$\pi^*(x, t) = H^{-1} B(x)^T \frac{\lambda \nabla_x \Psi(x, t)}{\Psi(x, t)} \quad (6)$$

As a result, obtaining Ψ is the most challenging computation in this form of problem.

If the optimal controls at certain time points, say $\{t_{1, \dots, n}\}$ with $t_n = T$, is emphasized, a representation in terms of the finite dimensional distribution $\Psi_i(x) = \Psi(x, t_i)$ can be obtained

from Eq. (4), only by computing the set $X = (X^0(t_0), \dots, X^0(t_n))$. The following recursive expression can be specifically obtained using the Markov property of $X^0(t)$ and marginalising intermediate states:

$$\Psi_i(x_{t_i}) = E_{X_{i+1}|x_{t_i}} \left[\Phi_i(x_{t_i}, X_{i+1}) \Psi_{i+1}(X_{i+1}) \right] \quad (7)$$

where,

$$\Phi_i(x_{t_i}, x_{t_{i+1}}) = E_{X^0(0)|x_{t_i}, x_{t_{i+1}}} \left[e^{-\frac{1}{\lambda} \int_{t_i}^{t_{i+1}} C(X^0(s), s) ds} \right] \quad (8)$$

with the expectation taken w.r.t. uncontrolled paths from x_{t_i} to $x_{t_{i+1}}$. It should be pointed out that $-\lambda \log \Phi_i$ can be considered as the (optimal)

expected cost for the problem ranging from x_{t_i} to $x_{t_{i+1}}$ over the time horizon $[t_i, t_{i+1}]$ under dynamics and running costs corresponding to those of the overall problem given in Eq. (2). Therefore, the problem naturally separated into simpler compounds; a set of short horizon problems – or indeed a nested hierarchy of such Φ – while, in the same time, it is a set of recursive evaluations backwards in time.

3 The Artificial Intelligence Embedding

Here, a demonstration is provided about the possibility of expressing Eq. (7) in terms of linear operators in RKHS. More details about the theory of RKHS and basic concepts of the presented theory in this paper can be found in [69-73] and [74-78] and [79-90], respectively.

Initially, the evaluation of a single step, i.e., Ψ_i given Ψ_{i+1} , is considered and model-based analytical expressions for the evaluation of Ψ_i in terms of certain operators in RKHS is extracted.

3.1 Embedding of the Artificial Intelligence by a Model-Based Analytical One Step Process

It shows that Eq. (7) may write as an inner product in a RKHS in the case based on the model. The process consists of three steps as expressing expectations in terms of inner products in a RKHS, firstly; adapting the basic expression to conditional expectations, secondly; and then, taking into account the conditional expectations of functions of both the conditional and conditioning variable.

H^κ is generally used as the RKHS of functions $Z \rightarrow R$ associated with the positive semi-definite kernel $\kappa(0,0)$. Considering P^Z as the set of random variables on Z , the embedding operator can be defined as $\varepsilon^\kappa : P^Z \rightarrow H^\kappa$:

$$E_Z[h(Z)] = \left\langle h, \underbrace{E_Z[k(Z,0)]}_{:=\varepsilon^\kappa[Z]} \right\rangle$$

$$\forall Z \in P^Z, h \in H^\kappa \quad (9)$$

in which the standard embedding of individual elements $z \in Z$ is directly extended into H^κ , given by $\varepsilon^\kappa[z] = \kappa(z,0)$ commonly encountered.

The main objective of the considered problem in this study is evaluating Ψ_i given in Eq. (7); i.e. proper embedding of $X_{i+1}|x_i$ which makes expressing the required expectation as an inner product in some RKHS possible. As $X_{i+1}|x_i$ is a simple random variable for fixed x_i , direct applying of Eq. (9) is possible. However, considering a general conditional random variable $Z|y$, act as a map $Y \rightarrow P^Z$ is more suitable. It yields random variables over Z which gives a value $y \in Y$, and defines a conditional embedding $U^{l\kappa} : H^l \rightarrow H^\kappa$ s.t.:

$$\varepsilon^\kappa[Z|y] = U^{l\kappa} \circ \varepsilon^l[y] \quad (10)$$

It should be pointed out that $\varepsilon^l[y] = l(0,y)$. It is used in kernel methods as the standard embedding operator of elements $y \in Y$ used in kernel methods. Therefore, conditional expectations can be expressed as:

$$E_{Z|y}[h(Z)] = \left\langle h, \underbrace{E_{Z|y}[\kappa(Z,0)]}_{:=\varepsilon^\kappa[Z|y]=U^{l\kappa}[l(y,0)]} \right\rangle \quad (11)$$

were demonstrated the existence and a specific form of an operator U which satisfies Eq. (10). It should be noted that the argument of the expectation in Eq. (7), and in particular of Φ , is a function of both the random variable, i.e., X_{i+1}

and the conditioning x_i . It is opposite to h in Eq. (11). Therefore, direct applying of Eq. (11) is not possible. However, an assisting random variable \tilde{X} is introduced where $P(\tilde{X}, X_{i+1}|x_i) = P(X_{i+1}|x_i)\delta_{\tilde{X}=x_i}$ with δ the delta distribution. Therefore, the following can be written for all $h \in H^\kappa$:

$$E_{X_{i+1}|x_i}[h(x_i, X_{i+1})] = E_{X_{i+1}, \tilde{X}|x_i}[h(\tilde{X}, X_{i+1})]$$

$$= \left\langle h, \varepsilon^\kappa[X_{i+1}, \tilde{X}|x_i] \right\rangle \quad (12)$$

If x_i is considered as a constant parameter, an alternative formulation can be obtained with equivalent analytical setting, but without an immediate yielding of a practical empirical estimator.

In order to substitute the specific argument encountered in Eq. (7) for the generic function h , it is assumed that H^ψ, H^ϕ , s.t. $\Psi \in H^\psi, \Phi \in H^\phi$, are given $(R^{d_x} \times R^{d_x} \rightarrow R$ and $R^{d_x} \rightarrow R$ are a space of functions H^ϕ and

functions H^ψ , respectively). The mismatch in the arity of functions in these spaces can be considered by extending H^ψ to $H^{\psi'}$; a space of functions $R^{d_x} \times R^{d_x} \rightarrow R$, using the kernel $\psi'((u,v), (u',v')) = \psi(u,u')$. It means that H^ψ and its tensor product with the RKHS of constant functions are identified. Then, Eq. (7) can be rewritten by taking the embedding of $X_{i+1}, \tilde{X} |_{x_i}$ into $H^\omega = H^\phi \otimes H^{\psi'}$ in which the product function of Φ_i, Ψ_{i+1} locates, and using Eqs. (11) and (12):

$$\begin{aligned} \Psi_i(x) &= E_{X_{i+1}|X_i} = x \left[\Phi_i(X_{i+1}, x) \cdot \Psi_{i+1}(X_{i+1}) \right] \\ &= \left\langle \Phi_i \otimes \Psi_{i+1}, \varepsilon^\omega [X_{i+1}, \tilde{X} |_{X_i = x}] \right\rangle \\ &= \left\langle \Phi_i \otimes \Psi_{i+1}, U^{\omega\kappa} \circ \varepsilon^\kappa [x] \right\rangle \end{aligned} \quad (13)$$

where κ is some kernel over chosen R^{d_x} . Taking κ so that it be able to reuse the pre-computed matrices during recursive evaluation of Ψ estimates of Ψ is computationally critical (see Eq. (16)).

3.2 Estimations of Model-Free Finite Sample

To evaluate the embedding of random variables U and artificial intelligence Eq. (13), it is necessary that expectations of kernels to be evaluated and remained inflexible. At the same time, there should be a detailed analytical knowledge about the system dynamics Eq. (1). However, forming empirical estimates is simple and it leads to practical algorithms.

Considering $D = \{(x, x')_{1..m}\}$ as a set of i.i.d, transition samples of the uncontrolled dynamics, e.g., a sample set obtained from trajectory executions under the policy π^0 , can be expressed. [80-90] were shown that a regularized estimate of $U^{\kappa\omega}$ can be defined as:

$$\hat{U}^{\psi\omega} = g_D^\omega (G_{\chi\chi}^\psi + \varepsilon m I)^{-1} g_\chi^\psi \quad (14)$$

where ε is a regularization parameter and $g_\chi^\psi, g_{\chi'}^\omega$ and $G_{\chi\chi}^\psi$ are the vectors of embeddings and Gramian on the sample data D respectively, i.e., $[g_\chi^\psi]_i = \phi(x_i, 0)$ and $[G_{\chi\chi}^\psi]_{ij} = \psi(x_i, x_j)$. The representations of Φ_i and Ψ_{i+1} in their respective RKHSs are necessary for evaluating Eq. (13). It can be assumed, due to recursive evaluation Ψ , that the empirical estimate of Ψ_{i+1} is defined as $\bar{\Psi}_{i+1} = \sum_{x_j \in \chi} [\alpha_i + 1] \psi(x_j, 0) = g_\chi^\psi \alpha_i$

where α_{i+1} is a vector of weights. In the same way, the representation of Φ_i in H^ϕ will be assumed as $g_B^\phi \beta$ for some weights β and B will be set. It will be shown that despite assurance from existence of such a representation with assuming that $\Phi_i \in H^\phi$, explicit computing of it is not necessary. By substituting the empirical operator $\hat{U}^{\kappa\omega}$ of Eq. (14) and the kernel expansions of Φ_i and Ψ_{i+1} into Eq. (13), matrix algebra yields the empirical estimate of $\Psi_i(x)$ as:

$$\bar{\Psi}_i(x) = \langle \psi(x, 0), g_\chi^\psi \alpha_i \rangle = G_{x\chi}^\psi \alpha_i \quad (15)$$

with weights α_i given by:

$$\alpha_i = \left[\begin{array}{c} G_{DB}^\phi \beta \odot G_{\chi'A}^\psi \alpha_{i+1} \\ = \Phi_i(\chi, \chi') \quad = \bar{\Psi}_{i+1}(\chi') \end{array} \right]^T (G_{\chi\chi}^\psi + \varepsilon m I)^{-1} \quad (16)$$

where \odot , denotes the Hadamard product. In addition to obtaining corresponding representation of Φ_i given by β , Eq. (16) can be used to recursive evaluation of the weights α . It should be noted that the term involving the representation of Φ_i can be written as:

$$G_{DB}^\phi \beta = \Phi_i(\mathcal{X}, \mathcal{X}') = (\Phi_i(x_1, x'_1), \Phi_i(x_1, x'_2), \dots)^T \quad (17)$$

This representation states that it is not necessary to obtain an explicit representation of Φ in terms of the kernel expansion in H^ϕ but the ability to evaluate Φ_i at the sample data D , which is comparable to evaluating the cost function, is sufficient.

It should be critically noted that direct recursive computation of all $\bar{\Psi}_{1..n}$ is possible since $\bar{\Psi}_i$ is a finite weighted sum of kernels and therefore, $\bar{\Psi}_i \in H^\psi$. In addition, pre-computing of all required matrices is possible since those are only functions of the sample data. Hence, it is easy to obtain an approximate optimal policy from Eq. (5) for fine discretisations of the problem.

4 Efficient Estimators

High computational complexity of $O(m^3)$ for the matrix inversion is one of main shortcomings of the basic estimator Eq. (16). If the same D and $O(m^2)$ per iteration are used in each time step, it is required. Another deficiency of Eq. (16) is that it needs to sample data under the uncontrolled dynamics. Therefore, off policy learning is not allowed. Alternative estimators for U by low rank approximations or importance sampling are used to remove these drawbacks. We choose to omit the discussion of these in order to address a, in our opinion, often overlooked aspect of efficiency when solving varying problems under the same dynamics. Actually, there are not individual tasks to be solved but several related problems should be solved in a repeated manner. For example, use of an optimized single reaching movement is limited due to need for a series of such movements with changing start and target states as a result of complicated interactions. The solution for this

problem in previous methods is generally based on the re-initialization for each specified problem. For instance, the start state of Monte-Carlo method is a new sample set, even when trivial changes are performed. Some extensions to the proposed method which aims to improve the efficiency of sampling as samples can reuse over and over for repeated applications are discussed in the following section.

4.1 Using Transition Sample Re-Use for Transfer Learning

One of the practical problems of estimators is the necessity of evaluating Φ at the training transitions. Hence, the favorite is obtaining an estimator based on evaluation of Φ on a separate, ideally arbitrary, data set D' . It can be seen that:

$$G_{DB}^\phi \beta = \langle \Phi, \phi(D, 0) \rangle = \langle \Phi, C_{ZZ}^{\phi\phi} (C_{ZZ}^{\phi\phi})^{-1} \phi(D, 0) \rangle \\ \approx \underbrace{\beta^T G_{BD'}}_{\Phi(D')} (G_{DD'}^\phi + \epsilon m I)^{-1} G_{DD}^\phi$$

where Z is some free random variable with support on $R^{d_x} \times R^{d_x}$, which is implemented as an empirical estimator based on a data set $D' = \{(x, x')_{1..m'}\}$ of i.i.d. samples from Z

(often in practice $D' \subseteq D$). The considered result can be simply achieved by substituting into Eq. (16) since the indicated evaluation of the r.h.s. only needs evaluation of D' at elements of Φ . Specifically, the ability of pre-computing and reusing the inverse matrix of Eq. (16) across changing tasks is obtained in addition to across different time steps by an assumed time stationary dynamics. This is of importance for efficient estimation in, e.g., the RL setting where incurred costs are known only at observed transitions or in cases where Φ can be freely evaluated but it is expensive to do so. However, generating large sets of transition samples may be comparatively cheap, e.g., the case of simple kinematic control where

cost evaluation requires collision detection. Note that this form makes explicit use of the kernel ϕ , and while we may not be able to guarantee $\Phi \in H^\phi$, by choosing a kernel such that the projection of Φ onto H^ϕ is close to Φ , we can expect good results.

4.2 Sampling by Augmenting Task

Generally, samples should be collected from the task agnostic dynamics X^0 . However, a task often induces regularities which suggest more suitable sampling distributions. To concentrate samples in regions of high Φ , considering the role Φ takes in Eq. (16) (via Eq. (17)) as a weight vector appears to be advantageous similar to importance sampling. It is clear that Φ is a suitable guidance for choosing the sampling distribution. However, for repeated task, Φ can be incorporated, partly, into the sampling process which allows incremental learning of the task.

For executing several tasks of a generic skill, it is frequently characterized by two components as an invariant cost relating to the skill and a task specific cost components. Assuming that the state cost decomposes as:

$$C(x, \theta, t) = C_{skill}(x, t) + C_{task}(x, t, \theta) \quad (18)$$

where θ parametrises the task. In this case, the path integral Eq. (4) can be rewritten as:

$$\Psi = E_{X^v(0)|x_t} \left[e^{-\int_t^T \frac{1}{\lambda} C_{task}(X^v(t), \theta, t)} \Psi(X^v(T), T) \right] \quad (19)$$

where, here, the expectation is taken w.r.t. path of X^v , which are the dynamics under the optimal policy under the invariant skill component of the cost. As a result of this, both can incrementally learned and using the previous results, the transfer of samples between varying tasks sharing a skill component.

5 Experimental Verification

5.1 Double Slit

To show Monte-Carlo approaches to artificial intelligence control, the double slit problem which is previously investigated is considered since it is, in one hand, is so simple to allow for a closed form solution for Ψ to be obtained, but on the other hand, it is so complicated to underline the deficiencies of some previous approaches. It deals with a particle moving with constant velocity in one coordinate, and simultaneously, its position in an orthogonal direction is influenced by noise and controls. The goal of the task is that square error should be reduced to a target value at the end time whereas avoiding obstacles at some intermediate time. For this problem, the one dimensional dynamics are $dx = u + d\zeta$ and the cost is given by:

$$C_\bullet(x) = \omega(x - x_{target})^2$$

$$C(x, t) = \begin{cases} 10^4 & \text{if } t = \frac{T}{2} \text{ and } x \in \text{Obstacle} \\ 0 & \text{else} \end{cases}$$

where ω is a weight. In this regard, a discretisation with time step $0.02s$ is considered.

A comparison is made between the true optimal policy and those obtained using two variants of the proposed estimator, as $\bar{\Psi}_{OC}$, $\bar{\Psi}_{RL}$ which are based on single transitions from uniformly sampled start states and a RL setting, learning from trajectory data without access to the cost, respectively. However, these used knowledge of the cost function to evaluate Φ in each step and the approach for sample sharing across time steps discussed, respectively. The low rank approximation and square exponential kernels $\psi(x, y) = \exp\left\{-\frac{(x - y)^2}{\gamma}\right\}$ with γ , equal to the mean distance of the data are used for both cases. Moreover, two alternative approaches are considered for more comparison. The first alternative is the trajectory based Monte-Carlo approach. While it has the same number of trajectories as used in the RL setting, it uses a

Laplace approximation to the true Ψ to attain a linear approximation of the optimal policy. The good performance of the proposed method can be clearly seen in which policies are considerably enhanced compared to those obtained by the alternative Monte-Carlo approach. However, the results of the proposed method are comparable to those obtained from the Laplace approximation. It should be pointed out that the results based on the Laplace approximation are computed by a prior knowledge about the true Ψ . Furthermore, it can be observed that the proposed approach makes better use of the sample provided by finding a policy which is applicable for varying starting positions. However, the multi-modality of the optimal policy is not recognized by the Monte-Carlo approach as improper results are obtained. For the variational approximation also the re-compute for each new starting location is very crucial.

The dependency of the estimate on the sample size is studied by comparing the evolution of the L_1 error of the estimates of Ψ at time $t = 0$. Sample size is denoted as total number of transition samples seen. Therefore, the number of trajectories is the sample size divided by 100 for $\bar{\Psi}^{RL}$. In order to also highlight the advantages of the sample re-use afforded by the approach in current paper, we also compare with $\bar{\Psi}$, the basic estimator given data of the same form as $\bar{\Psi}^{RL}$, i.e. recursive application of Eq. (16) without sample sharing across time steps.

5.2 Reaching Task of Arm Subspace

For simulating constrained tasks, reaching tasks on a subspace of the end-effector space of a torque controlled 5dof arm is considered. The considered skill component includes moving with the end-effector close to a two dimensional task space, while the task examples are given by specific reach targets. A linear subspace of the end effector space which is a non-linear subspace of the joint space is used in this section and the cost consists of two components:

$$C_{skill}(x, t) = \omega_{skill} \|J_{\varphi}(x) - j\|^2$$

$$C_{task}(x, \theta) = \omega_{task} \|\varphi(x) - \theta\|^2$$

where $\varphi(0)$ is the mapping from joint to end-effector coordinates, J & j define the task subspace, θ specifies the reaching target and ω 's are weights. A position control over a 2s horizon with a 0.02s discretisation is considered here, again.

Due to the restrictions of low cost trajectories to a small subspace, this task is challenging for sample based approaches. Although some researchers were suggested to improve end-effector exploration an inverse dynamics policy can be used, it is necessary to increase sample sizes for overcoming this problem. With focus on the case of changing targets, the ideas originated are used as it is assumed that the operators have been estimated under the skill augmented dynamics (as an alternative to explicit sample generating, utilizing the importance sample based estimator and then, collecting a sample under X^0 , which is more time consuming than what that is performed here), and then, the subsequent learning for a novel task is considered by the estimator, utilizing the already estimated operators in two ways. They are directly used in the calculation of $\bar{\Psi}$; however, it should be noted that the trajectories can be sampled without considering a specific policy since these are only required to provide D' . As a result, the policy arising is used only when considering C_{skill} , i.e., the skill policy associated with $\bar{\Psi}$ computed using the given operators and $C_{task}(0) = 0$.

It is demonstrated that sampling under the skill policy is more effective in exploring the task relevant sub space than null policy.

6 Summary and Conclusion

An innovative approach is presented in the current paper to solve stochastic optimal control problems which have the artificial intelligence control form

using Monte-Carlo estimates of operators arising from a RKHS embedding of the problem. It leads to a consistent estimate of problem Ψ . Although direct application of Monte-Carlo estimation to point evaluation of Ψ also obtains a consistent estimate, a trajectory sample for each state at which an action is to be computed is needed due to impracticality of computing the controls for anything but simple problems. Despite some suggestions provided by previous works to reduce sample complexity, it is demonstrated that the proposed approach is of more generality in policy than these previous suggestions which are not consistent in their processes. Moreover, it is shown that sample re-using is possible by the proposed estimators, particularly for cases in which a new sample set is needed in advance. Particular emphasis is focused on the transfer in cases where execution of several, potentially related, tasks on the same plant is required. It shows that use of samples from all tasks to learn invariant aspects are possible. An alternative approach of the proposed method can be defined to combine solutions to local control problem, as defined by Φ , and hence, to solve a more complicated large scale problem. Combining the proposed methods with alternative variation approaches will be followed in future studies so that a good estimate for the comparatively simpler local problems can be achieved. Although kernel is not chosen in this study, making informed kernel choices based on in advance knowledge about the structure of the problem may be positive. In current work, a method of utilizing RKHS embeddings of the transition probability in computing the value functions in MDPs. However, this work is of some advantages over their work. The first advantage is that the optimal controls are directly obtained in this paper instead of computing the value function. It is better than use of explicit maximization to obtain optimal controls. In addition, in current study focused on finite state problems (where computation of the optimal u is simpler), while harder continuous problem are studied in this paper which provides convergence guarantees in this setting. As the final point, it should be noted that the structure of the problem is

used to efficiently estimates the required quantities which leads to efficient sample re-use and transfer.

Acknowledgment

The authors would like to thank Dr. Taher Lotfi for help and valuable discussions.

References

- [1] G. E. Fasshauer, F. J. Hickernell, Q. Ye, Solving support vector machines in reproducing kernel Banach spaces with positive definite functions, *Applied and Computational Harmonic Analysis*, Volume 38, Issue 1, January 2015, Pages 115-139.
- [2] F.Z. Geng, S.P. Qian, M.G. Cui, Improved reproducing kernel method for singularly perturbed differential-difference equations with boundary layer behavior, *Applied Mathematics and Computation*, Volume 252, 1 February 2015, Pages 58-63.
- [3] M. L. L. García, R. García-Ródenas, A. G. Gómez, K-means algorithms for functional data, *Neurocomputing*, Volume 151, Part 1, 3 March 2015, Pages 231-245.
- [4] S. Lv, F. Zhou, Optimal learning rates of γ -type multiple kernel learning under general conditions, *Information Sciences*, Volume 294, 10 February 2015, Pages 255-268.
- [5] A. Chattopadhyay, B. K. Das, J. Sarkar, Tensor product of quotient Hilbert modules, *Journal of Mathematical Analysis and Applications*, Volume 424, Issue 1, 1 April 2015, Pages 727-747.
- [6] P.E.T. Jorgensen, R. Niedzialomski, Extension of positive definite functions, *Journal of Mathematical Analysis and Applications*, Volume 422, Issue 1, 1 February 2015, Pages 712-740.
- [7] J. Dick, P. Kritzer, G. Leobacher, F. Pillichshammer, A reduced fast component-by-component construction of lattice points for integration in weighted spaces with fast decreasing weights, *Journal of Computational and Applied Mathematics*, Volume 276, 1 March 2015, Pages 1-15.
- [8] T. Villmann, S. Haase, M. Kaden, Kernelized vector quantization in gradient-descent learning, *Neurocomputing*, Volume 147, 5 January 2015, Pages 83-95.
- [9] J. Ma, J. Zhao, Y. Ma, J. Tian, Non-rigid visible and infrared face registration via regularized Gaussian fields criterion, *Pattern Recognition*, Volume 48, Issue 3, March 2015, Pages 772-784.
- [10] H. Shin, S. Lee, Canonical correlation analysis for irregularly and sparsely observed functional data, *Journal of*

Multivariate Analysis, Volume 134, February 2015, Pages 1-18.

[11] E. Novak, M. Ullrich, H. Woźniakowski, Complexity of oscillatory integration for univariate Sobolev spaces, *Journal of Complexity*, Volume 31, Issue 1, February 2015, Pages 15-41.

[12] M. Lange, M. Biehl, T. Villmann, Non-Euclidean principal component analysis by Hebbian learning, *Neurocomputing*, Volume 147, 5 January 2015, Pages 107-119.

[13] Y. Hao, A. Wang, L. Zhang, On holomorphic isometric immersions of nonhomogeneous Kähler–Einstein manifolds into the infinite dimensional complex projective space, *Journal of Mathematical Analysis and Applications*, Volume 423, Issue 1, 1 March 2015, Pages 547-560.

[14] L. Li, Y. Nakada, D. Nestor, W. Shang, R. Weir, Normal weighted composition operators on weighted Dirichlet spaces, *Journal of Mathematical Analysis and Applications*, Volume 423, Issue 1, 1 March 2015, Pages 758-769.

[15] X. Li, H. He, Z. Yin, F. Chen, J. Cheng, KPLS-based image super-resolution using clustering and weighted boosting, *Neurocomputing*, Volume 149, Part B, 3 February 2015, Pages 940-948.

[16] D. Alpay, F. Colombo, D. P. Kimsey, I. Sabadini, An extension of Herglotz's theorem to the quaternions, *Journal of Mathematical Analysis and Applications*, Volume 421, Issue 1, 1 January 2015, Pages 754-778.

[17] D. Grob, R.S. Kraußhar, A Selberg trace formula for hypercomplex analytic cusp forms, *Journal of Number Theory*, Volume 148, March 2015, Pages 398-428.

[18] X. Zhan, D. Ghosh, Incorporating auxiliary information for improved prediction using combination of kernel machines, *Statistical Methodology*, Volume 22, January 2015, Pages 47-57.

[19] X. Tian, L. Jiao, L. Yi, K. Guo, X. Zhang, The image segmentation based on optimized spatial feature of superpixel, *Journal of Visual Communication and Image Representation*, Volume 26, January 2015, Pages 146-160.

[20] R.V. Bessonov, Duality theorems for coinvariant subspaces of, *Advances in Mathematics*, Volume 271, 5 February 2015, Pages 62-90.

[21] J. Chen, W. Pedrycz, M. Ha, L. Ma, Set-valued samples based support vector regression and its applications, *Expert Systems with Applications*, Volume 42, Issue 5, 1 April 2015, Pages 2502-2509.

[22] Z. Chen, W. Zuo, Q. Hu, L. Lin, Kernel sparse representation for time series classification, *Information Sciences*, Volume 292, 20 January 2015, Pages 15-26.

[23] U. M. Al-Saggaf, M. Moinuddin, M. Arif, A. Zerguine, The q-Least Mean Squares algorithm, *Signal Processing*, Volume 111, June 2015, Pages 50-60.

[24] N. H. Tuan, T. T. Binh, N. D. Minh, T. T. Nghia, An improved regularization method for initial inverse problem in 2-D heat equation, *Applied Mathematical Modelling*, Volume 39, Issue 2, 15 January 2015, Pages 425-437.

[25] K. Slavakis, P. Bouboulis, S. Theodoridis, Chapter 17 - Online Learning in Reproducing Kernel Hilbert Spaces, In: Paulo S.R. Diniz, Johan A.K. Suykens, R. Chellappa and S. Theodoridis, Editor(s), *Academic Press Library in Signal Processing*, Elsevier, 2014, Volume 1, Pages 883-987.

[26] T. Jordão, V.A. Menegatto, Weighted Fourier–Laplace transforms in reproducing kernel Hilbert spaces on the sphere, *Journal of Mathematical Analysis and Applications*, Volume 411, Issue 2, 15 March 2014, Pages 732-741.

[27] J. González, I. Vujačić, E. Wit, Reproducing kernel Hilbert space based estimation of systems of ordinary differential equations, *Pattern Recognition Letters*, Volume 45, 1 August 2014, Pages 26-32.

[28] M. Seto, S. Suda, T. Taniguchi, Gram matrices of reproducing kernel Hilbert spaces over graphs, *Linear Algebra and its Applications*, Volume 445, 15 March 2014, Pages 56-68.

[29] H. Woracek, Reproducing kernel almost Pontryagin spaces, *Linear Algebra and its Applications*, Volume 461, 15 November 2014, Pages 271-317.

[30] X. Luo, Z. Lu, X. Xu, Reproducing kernel technique for high dimensional model representations (HDMMR), *Computer Physics Communications*, Volume 185, Issue 12, December 2014, Pages 3099-3108.

[31] H. Du, G. Zhao, C. Zhao, Reproducing kernel method for solving Fredholm integro-differential equations with weakly singularity, *Journal of Computational and Applied Mathematics*, Volume 255, 1 January 2014, Pages 122-132.

[32] F.Z. Geng, S.P. Qian, A new reproducing kernel method for linear nonlocal boundary value problems, *Applied Mathematics and Computation*, Volume 248, 1 December 2014, Pages 421-425.

[33] M. Gnewuch, S. Mayer, K. Ritter, On weighted Hilbert spaces and integration of functions of infinitely many variables, *Journal of Complexity*, Volume 30, Issue 2, April 2014, Pages 29-47.

- [34] G. Pillonetto, F. Dinuzzo, T. Chen, G. De Nicolao, L. Ljung, Kernel methods in system identification, machine learning and function estimation: A survey, *Automatica*, Volume 50, Issue 3, March 2014, Pages 657-682.
- [35] M. Mitkovski, B. D. Wick, A reproducing kernel thesis for operators on Bergman-type function spaces, *Journal of Functional Analysis*, Volume 267, Issue 7, 1 October 2014, Pages 2028-2055.
- [36] O. Abu Arqub, M. Al-Smadi, Numerical algorithm for solving two-point, second-order periodic boundary value problems for mixed integro-differential equations, *Applied Mathematics and Computation*, Volume 243, 15 September 2014, Pages 911-922.
- [37] E. De Vito, L. Rosasco, A. Toigo, Learning sets with separating kernels, *Applied and Computational Harmonic Analysis*, Volume 37, Issue 2, September 2014, Pages 185-217.
- [38] J. Shawe-Taylor, S. Sun, Chapter 16 - Kernel Methods and Support Vector Machines, In: Paulo S.R. Diniz, Johan A.K. Suykens, R. Chellappa and S. Theodoridis, Editor(s), *Academic Press Library in Signal Processing*, Elsevier, 2014, Volume 1, Pages 857-881.
- [39] S. H. Kang, B. Shafei, G. Steidl, Supervised and transductive multi-class segmentation using p-Laplacians and RKHS methods, *Journal of Visual Communication and Image Representation*, Volume 25, Issue 5, July 2014, Pages 1136-1148.
- [40] J. Xu, Y. Y. Tang, B. Zou, Z. Xu, L. Li, Y. Lu, Generalization performance of Gaussian kernels SVMC based on Markov sampling, *Neural Networks*, Volume 53, May 2014, Pages 40-51.
- [41] S. Momani, O. Abu Arqub, T. Hayat, H. Al-Sulami, A computational method for solving periodic boundary value problems for integro-differential equations of Fredholm-Volterra type, *Applied Mathematics and Computation*, Volume 240, 1 August 2014, Pages 229-239.
- [42] H. Fan, Q. Song, S. B. Shrestha, Online learning with kernel regularized least mean square algorithms, *Knowledge-Based Systems*, Volume 59, March 2014, Pages 21-32.
- [43] Y. Fu, L. Li, U. Kaehler, P. Cerejeiras, On the Fock space of metaanalytic functions, *Journal of Mathematical Analysis and Applications*, Volume 414, Issue 1, 1 June 2014, Pages 176-187.
- [44] H. Xue, S. Chen, Discriminality-driven regularization framework for indefinite kernel machine, *Neurocomputing*, Volume 133, 10 June 2014, Pages 209-221.
- [45] T. J. Hansen, T. J. Abrahamsen, L. K. Hansen, Denoising by semi-supervised kernel PCA preimaging, *Pattern Recognition Letters*, Volume 49, 1 November 2014, Pages 114-120.
- [46] G. Rozenblum, N. Vasilevski, Toeplitz operators defined by sesquilinear forms: Fock space case, *Journal of Functional Analysis*, Volume 267, Issue 11, 1 December 2014, Pages 4399-4430.
- [47] S. Li, C. Deng, W. Sun, The isometric identities and inversion formulas of complex continuous wavelet transforms, *Applied Mathematics and Computation*, Volume 233, 1 May 2014, Pages 116-126.
- [48] Y. Mo, T. Qian, Support vector machine adapted Tikhonov regularization method to solve Dirichlet problem, *Applied Mathematics and Computation*, Volume 245, 15 October 2014, Pages 509-519.
- [49] J. Dick, M. Gnewuch, Optimal randomized changing dimension algorithms for infinite-dimensional integration on function spaces with ANOVA-type decomposition, *Journal of Approximation Theory*, Volume 184, August 2014, Pages 111-145.
- [50] D.-R. Chen, H. Li, Convergence rates of learning algorithms by random projection, *Applied and Computational Harmonic Analysis*, Volume 37, Issue 1, July 2014, Pages 36-51.
- [51] R. G. Douglas, Y.-S. Kim, H.-K. Kwon, J. Sarkar, Curvature invariant and generalized canonical operator models – II, *Journal of Functional Analysis*, Volume 266, Issue 4, 15 February 2014, Pages 2486-2502.
- [52] J. Peng, L. Li, Support vector regression in sum space for multivariate calibration, *Chemometrics and Intelligent Laboratory Systems*, Volume 130, 15 January 2014, Pages 14-19.
- [53] J. A. Ball, D. S. Kaliuzhnyi-Verbovetskyi, Rational Cayley inner Herglotz-Agler functions: Positive-kernel decompositions and transfer-function realizations, *Linear Algebra and its Applications*, Volume 456, 1 September 2014, Pages 138-156.
- [54] X. Yang, A. Cao, Q. Song, G. Schaefer, Y. Su, Vicinal support vector classifier using supervised kernel-based clustering, *Artificial Intelligence in Medicine*, Volume 60, Issue 3, March 2014, Pages 189-196.
- [55] V.A. Menegatto, Differentiability of bizonal positive definite kernels on complex spheres, *Journal of Mathematical Analysis and Applications*, Volume 412, Issue 1, 1 April 2014, Pages 189-199.
- [56] O. González-Recio, G.J.M. Rosa, D. Gianola, Machine learning methods and predictive ability metrics for genome-

wide prediction of complex traits, *Livestock Science*, Volume 166, August 2014, Pages 217-231.

[57] V. Kůrková, P. C. Kainen, Comparing fixed and variable-width Gaussian networks, *Neural Networks*, Volume 57, September 2014, Pages 23-28.

[58] J. L. Rojo-Álvarez, M. Martínez-Ramón, J. Muñoz-Marí, G. Camps-Valls, A unified SVM framework for signal estimation, *Digital Signal Processing*, Volume 26, March 2014, Pages 1-20.

[59] F.Z. Geng, S.P. Qian, S. Li, A numerical method for singularly perturbed turning point problems with an interior layer, *Journal of Computational and Applied Mathematics*, Volume 255, 1 January 2014, Pages 97-105.

[60] L. Wang, H. Shi, Improved Kernel PLS-based Fault Detection Approach for Nonlinear Chemical Processes, *Chinese Journal of Chemical Engineering*, Volume 22, Issue 6, June 2014, Pages 657-663.

[61] A. Haimi, Bulk asymptotics for polyanalytic correlation kernels, *Journal of Functional Analysis*, Volume 266, Issue 5, 1 March 2014, Pages 3083-3133.

[62] D. Scheinker, Hilbert function spaces and the Nevanlinna-Pick problem on the polydisc II, *Journal of Functional Analysis*, Volume 266, Issue 1, 1 January 2014, Pages 355-367.

[63] S. Huang, L. Jin, Y. Fang, X. Wei, Online heterogeneous feature fusion machines for visual recognition, *Neurocomputing*, Volume 123, 10 January 2014, Pages 100-109.

[64] A. Haimi, H. Hedenmalm, Asymptotic expansion of polyanalytic Bergman kernels, *Journal of Functional Analysis*, Volume 267, Issue 12, 15 December 2014, Pages 4667-4731.

[65] J. L. Godoy, D.A. Zumoffen, J. R. Vega, J. L. Marchetti, New contributions to non-linear process monitoring through kernel partial least squares, *Chemometrics and Intelligent Laboratory Systems*, Volume 135, 15 July 2014, Pages 76-89.

[66] K. T. Mynbaev, S. Nadarajah, C. S. Withers, A. S. Aipenova, Improving bias in kernel density estimation, *Statistics & Probability Letters*, Volume 94, November 2014, Pages 106-112.

[67] H. Chen, Z. Pan, L. Li, Learning performance of coefficient-based regularized ranking, *Neurocomputing*, Volume 133, 10 June 2014, Pages 54-62.

[68] M. Fan, N. Gu, H. Qiao, B. Zhang, Dimensionality reduction: An interpretation from manifold regularization perspective, *Information Sciences*, Volume 277, 1 September 2014, Pages 694-714.

[69] W. Wang, B.-G. Hu, Z.-F. Wang, Globality and locality incorporation in distance metric learning, *Neurocomputing*, Volume 129, 10 April 2014, Pages 185-198.

[70] C. Cortes, M. Mohri, Domain adaptation and sample bias correction theory and algorithm for regression, *Theoretical Computer Science*, Volume 519, 30 January 2014, Pages 103-126.

[71] H. Jin, X. Chen, J. Yang, L. Wu, Adaptive soft sensor modeling framework based on just-in-time learning and kernel partial least squares regression for nonlinear multiphase batch processes, *Computers & Chemical Engineering*, Volume 71, 4 December 2014, Pages 77-93.

[72] E. Carneiro, F. Littmann, Extremal functions in de Branges and Euclidean spaces, *Advances in Mathematics*, Volume 260, 1 August 2014, Pages 281-349.

[73] Z. Chen, S. Xiong, Z. Fang, Q. Li, B. Wang, Q. Zou, A kernel support vector machine-based feature selection approach for recognizing Flying Apsaras' streamers in the Dunhuang Grotto Murals, China, *Pattern Recognition Letters*, Volume 49, 1 November 2014, Pages 107-113.

[74] J.A.K. Suykens, Chapter 13 - Introduction to Machine Learning, In: Paulo S.R. Diniz, Johan A.K. Suykens, R. Chellappa and S. Theodoridis, Editor(s), *Academic Press Library in Signal Processing*, Elsevier, 2014, Volume 1, Pages 765-773.

[75] G. Gnecco, M. Gori, S. Melacci, M. Sanguineti, A theoretical framework for supervised learning from regions, *Neurocomputing*, Volume 129, 10 April 2014, Pages 25-32.

[76] J. A. Nichols, F. Y. Kuo, Fast CBC construction of randomly shifted lattice rules achieving convergence for unbounded integrands over in weighted spaces with POD weights, *Journal of Complexity*, Volume 30, Issue 4, August 2014, Pages 444-468.

[77] W.-J. Chen, Y.-H. Shao, N.-Y. Deng, Z.-L. Feng, Laplacian least squares twin support vector machine for semi-supervised classification, *Neurocomputing*, Volume 145, 5 December 2014, Pages 465-476.

[78] N. Singh, P. T. Fletcher, J. S. Preston, R. D. King, J.S. Marron, M. W. Weiner, S. Joshi, Alzheimer's Disease Neuroimaging Initiative (ADNI), Quantifying anatomical shape variations in neurological disorders, *Medical Image Analysis*, Volume 18, Issue 3, April 2014, Pages 616-633.

[79] A.E. Frazho, S. ter Horst, M.A. Kaashoek, State space formulae for stable rational matrix solutions of a Leech problem, *Indagationes Mathematicae*, Volume 25, Issue 2, 14 March 2014, Pages 250-274.

[80] A. Cabada, A.R. Hayotov, K.M. Shadimetov, Construction of -splines in space by Sobolev method,

Applied Mathematics and Computation, Volume 244, 1
October 2014, Pages 542-551.

[81] X.Y. Li, B.Y. Wu, A continuous method for nonlocal functional differential equations with delayed or advanced arguments, Journal of Mathematical Analysis and Applications, Volume 409, Issue 1, 1 January 2014, Pages 485-493.

[82] R. R. Coifman, M. J. Hirn, Diffusion maps for changing data, Applied and Computational Harmonic Analysis, Volume 36, Issue 1, January 2014, Pages 79-107.

[83] M. M. Meerschaert, F. Sabzikar, Stochastic integration for tempered fractional Brownian motion, Stochastic Processes and their Applications, Volume 124, Issue 7, July 2014, Pages 2363-2387.

[84] J. C. Príncipe, B. Chen, L. G. S. Giraldo, Chapter 24 - Information Based Learning, In: Paulo S.R. Diniz, Johan A.K. Suykens, R. Chellappa and S. Theodoridis, Editor(s), Academic Press Library in Signal Processing, Elsevier, 2014, Volume 1, Pages 1379-1414.

[85] S. Asserda, A. Hichame, Pointwise estimate for the Bergman kernel of the weighted Bergman spaces with exponential type weights, Comptes Rendus Mathématique, Volume 352, Issue 1, January 2014, Pages 13-16.

[86] M. T. Jury, Clark theory in the Drury–Arveson space, Journal of Functional Analysis, Volume 266, Issue 6, 15 March 2014, Pages 3855-3893.

[87] M.-H. Hsu, L.-C. Wang, Z. He, Interpolation problems for holomorphic functions, Linear Algebra and its Applications, Volume 452, 1 July 2014, Pages 270-280.

[88] J. Cui, Y. Duan, Berger measure for, Journal of Mathematical Analysis and Applications, Volume 413, Issue 1, 1 May 2014, Pages 202-211.

[89] X. Zhou, M. Belkin, Chapter 22 - Semi-Supervised Learning, In: Paulo S.R. Diniz, Johan A.K. Suykens, R. Chellappa and S. Theodoridis, Editor(s), Academic Press Library in Signal Processing, Elsevier, 2014, Volume 1, Pages 1239-1269.

[90] W. Xie, Z. Lu, Y. Peng, J. Xiao, Graph-based multimodal semi-supervised image classification, Neurocomputing, Volume 138, 22 August 2014, Pages 167-179.