



Hilbert Space Embeddings of POMDPs

Yu Nishiyama¹ Abdeslam Boularias² Arthur Gretton^{2,3} Kenji Fukumizu¹

¹The Institute of Statistical Mathematics

²Max Planck Institute for Intelligent Systems

³Gatsby Computational Neuroscience Unit, CSML, UCL



MAX-PLANCK-GESELLSCHAFT MAX-PLANCK-GESELLSCHAFT

1 Overview

A nonparametric approach for policy learning for POMDPs is proposed. The approach represents distributions over the states, observations, and actions as embeddings in feature spaces, which are reproducing kernel Hilbert spaces. Distributions over states given the observations are obtained by applying the kernel Bayes' rule to these distribution embeddings. Policies and value functions are defined on the feature space over states, which leads to a feature space expression for the Bellman equation. Value iteration may then be used to estimate the optimal value function and associated policy. Experimental results confirm that the correct policy is learned using the feature space representation.

2 Background

2.1 Partially Observable Markov Decision Process (POMDP)

A POMDP is a tuple $\langle \mathcal{S}, \mathcal{A}, T, R, \mathcal{O}, Z \rangle$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, T is the transition function with $T(s, a, s') = \Pr(s'|s, a)$ for $s', s \in \mathcal{S}$ and $a \in \mathcal{A}$, R is the reward function with $R(s, a)$ for $s \in \mathcal{S}$ and $a \in \mathcal{A}$, \mathcal{O} is the set of observations, and Z is the observation function with $Z(s, o) = \Pr(o|s)$ for $o \in \mathcal{O}$ and $s \in \mathcal{S}$.

2.2 Reinforcement Learning in POMDP Environments

Belief Update Rule:

Belief b is a distribution over \mathcal{S} . Given an initial belief b_0 and a history of actions and observations $h_{t+1} = \{a_0, o_1, \dots, a_t, o_{t+1}\}$, belief $b_{t+1}(t \geq 0)$ is updated according to Bayes' rule

$$b_{t+1}(s_{t+1}) = \frac{Z(s_{t+1}, o_{t+1})P(s_{t+1}|a_t; b_t)}{P(o_{t+1}|a_t; b_t)}, \quad (1)$$

where $P(s_{t+1}|a_t; b_t) = \mathbb{E}_{S_t \sim b_t} [T(S_t, a_t, s_{t+1})]$ and $P(o_{t+1}|a_t; b_t) = \mathbb{E}_{S_{t+1} \sim P(\cdot|b_t, a_t)} [Z(o_{t+1}, S_{t+1})]$. In what follows, $b_{t+1}(s_{t+1}) = b^{a, o'}$ when o' is observed after executing action a in belief b .

Goal:

A goal is to find an optimal policy $\pi^* : b \mapsto a$ that maximizes the expected sum of discounted rewards with infinite horizon $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R_t]$, where $\gamma \in (0, 1)$ is a discount factor.

Bellman Optimality Equation:

The optimal policy π^* and the value function $V^*(b)$ are the fixed point of the Bellman optimality equation

$$V^*(b) = \max_{a \in \mathcal{A}} Q^*(b, a), \quad \pi^*(b) = \arg \max_{a \in \mathcal{A}} Q^*(b, a), \quad Q^*(b, a) = \mathbb{E}_{S \sim b} [R(S, a)] + \gamma \mathbb{E}_{O' \sim P(\cdot|b, a)} [V^*(b^{a, O'})],$$

Value Iteration Algorithm:

The optimal policy π^* and the value function $V^*(b)$ are estimated by the value iteration algorithm $V_d = HV_{d-1}$ ($d \geq 1$), where H is the Bellman operator and V_d is the d -step value function. H is isotonic and contractive.

Initial Value:

Initial value $V_0(b)$ is set with initial Q -value $Q_0(s, a)$ as $V_0(b) = \max_{a \in \mathcal{A}} \mathbb{E}_{S \sim b(\cdot)} [Q_0(S, a)]$. Examples are reward $Q_0(s, a) = R(s, a)$ and QMDP approximation $Q_0(s, a) = Q^{MDP}(s, a)$ [?], where $Q^{MDP}(s, a)$ is the result of MDP value iteration, approximating the POMDP by an MDP.

2.3 Kernel Method for Probabilities

Mean Embedding:

The mean embedding of distribution P in $\mathcal{H}_{\mathcal{X}}$ is the mean features of P , i.e., the RKHS element $\mu_{\mathcal{X}} = \mathbb{E}_{X \sim P} [k_{\mathcal{X}}(X, \cdot)]$. $\langle \mu_{\mathcal{X}}, f \rangle_{\mathcal{H}_{\mathcal{X}}} = \mathbb{E}_{X \sim P} [f(X)]$ holds for all $f \in \mathcal{H}_{\mathcal{X}}$. Empirical estimate $\hat{\mu}_{\mathcal{X}}$ has a form $\hat{\mu}_{\mathcal{X}} = \Upsilon \alpha$, where $\Upsilon = (k_{\mathcal{X}}(\cdot, X_1), \dots, k_{\mathcal{X}}(\cdot, X_n))$ is a feature matrix and $\alpha = (\alpha_1, \dots, \alpha_n)^{\top} \in \mathbb{R}^n$ is a weight vector on samples (X_1, \dots, X_n) . $\mathbb{E}_{X \sim P} [f(X)]$ can then be nonparametrically estimated by $\langle \hat{\mu}_{\mathcal{X}}, f \rangle_{\mathcal{H}_{\mathcal{X}}} = \alpha^{\top} \mathbf{f}$, where $\mathbf{f} = (f(X_1), \dots, f(X_n))^{\top}$ is the sample vector of f .

Conditional Embedding Operators & Kernel Bayes' Rule (KBR):

Let $\mathcal{H}_{\mathcal{X}}$ and $\mathcal{H}_{\mathcal{Y}}$ be RKHSs associated with $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ over $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ and $(\mathcal{Y}, \mathcal{B}_{\mathcal{Y}})$, respectively. Let (X, Y) be a random variable taking values on $\mathcal{X} \times \mathcal{Y}$ with distribution P and the density $p(x, y)$. The conditional density functions $\{p(Y|X=x)|x \in \mathcal{X}\}$ define a family of embeddings $\{\mu_{Y|x}\}$ in $\mathcal{H}_{\mathcal{Y}}$. A mapping from $k_{\mathcal{X}}(x, \cdot) \in \mathcal{H}_{\mathcal{X}}$ to $\mu_{Y|x} \in \mathcal{H}_{\mathcal{Y}}$ for all $x \in \mathcal{X}$ can be characterized by conditional embedding operator $\mathcal{U}_{Y|\mathcal{X}} : \mathcal{H}_{\mathcal{X}} \rightarrow \mathcal{H}_{\mathcal{Y}}$,

$$\mu_{Y|x} = \mathcal{U}_{Y|\mathcal{X}} k_{\mathcal{X}}(x, \cdot) = C_{Y\mathcal{X}} C_{X\mathcal{X}}^{-1} k_{\mathcal{X}}(x, \cdot), \quad (2)$$

where $C_{Y\mathcal{X}}$ and $C_{X\mathcal{X}}$ are uncentred covariance operators with respect to P [?].

Since a posterior distribution is a conditional distribution, the embedding of a posterior can be expressed as a conditional embedding operator [?]. Let Π be a prior distribution with density $\pi(x)$, and (\bar{X}, \bar{Y}) be a new random variable with distribution Q and the density $q(x, y) = p(y|x)\pi(x)$. The embedding of a posterior $q(\bar{X}|\bar{Y}=y)$ given y is expressed by a conditional embedding operator $\mathcal{U}_{\bar{X}|\bar{Y}}$

$$\mu_{\bar{X}|y} = \mathcal{U}_{\bar{X}|\bar{Y}} k_{\bar{Y}}(y, \cdot) = C_{\bar{X}\bar{Y}} C_{\bar{Y}\bar{Y}}^{-1} k_{\bar{Y}}(y, \cdot), \quad (3)$$

where $C_{\bar{X}\bar{Y}}$ and $C_{\bar{Y}\bar{Y}}$ are covariance operators with respect to Q .

3 Kernel POMDP (kPOMDP)

3.1 Kernel Bellman Equations (KBes)

Let mean embeddings of relevant distributions $b(S)$, $P(S'|a; b)$, $P(O'|a; b)$, $b^{a, o'}(S')$ in the corresponding RKHSs $\mathcal{H}_{\mathcal{S}}$, $\mathcal{H}_{\mathcal{O}}$ be

$$\begin{aligned} \mu_{\mathcal{S}} &= \mathbb{E}_{S \sim b(\cdot)} [\varphi(S)], & \mu_{S'|a; \mu_{\mathcal{S}}} &= \mathbb{E}_{S' \sim P(\cdot|a; b)} [\varphi(S')], \\ \mu_{O'|a; \mu_{\mathcal{S}}} &= \mathbb{E}_{O' \sim P(\cdot|a; b)} [\phi(O')], & \mu_{S'}^{a, o'} &= \mathbb{E}_{S' \sim b^{a, o'}(\cdot)} [\varphi(S')]. \end{aligned}$$

Let $\mathcal{U}_{S'|S, \mathcal{A}}$ and $\mathcal{U}_{O|S}$ be conditional embedding operators for transition model T and observation model Z , respectively. $\mathcal{U}_{S'|O}^{(a, \mu_{\mathcal{S}})}$ be a posterior embedding operator with a prior embedding $\mu_{S'|a; \mu_{\mathcal{S}}}$ corresponding to eq.(3). These embedding operators yield relations

$$\mu_{S'|a; \mu_{\mathcal{S}}} = \mathcal{U}_{S'|S, \mathcal{A}} \mu_{\mathcal{S}} \otimes k_{\mathcal{A}}(a, \cdot), \quad \mu_{O'|a; \mu_{\mathcal{S}}} = \mathcal{U}_{O|S} \mu_{S'|a; \mu_{\mathcal{S}}}, \quad \mu_{S'}^{a, o'} = \mathcal{U}_{S'|O}^{(a, \mu_{\mathcal{S}})} k_{O}(o', \cdot). \quad (4)$$

Let $\mathcal{P}_{\mathcal{S}}$ be the set of beliefs and $\mathcal{I}_{\mathcal{S}}$ be the set of embeddings of $\mathcal{P}_{\mathcal{S}}$ in $\mathcal{H}_{\mathcal{S}}$.

Claim 1. Let $R(\cdot, a) \in \mathcal{H}_{\mathcal{S}}$ and $V^*(\mu_{S'}^{a, (\cdot)}) \in \mathcal{H}_{\mathcal{O}}$ for all $a \in \mathcal{A}$ and $\mu_{\mathcal{S}} \in \mathcal{I}_{\mathcal{S}}$. The kernel Bellman optimality equations on RKHS $\mathcal{H}_{\mathcal{S}}$ may be

$$\begin{aligned} V^*(\mu_{\mathcal{S}}) &= \max_{a \in \mathcal{A}} Q^*(\mu_{\mathcal{S}}, a), \quad \pi^*(\mu_{\mathcal{S}}) = \arg \max_{a \in \mathcal{A}} Q^*(\mu_{\mathcal{S}}, a). \\ Q^*(\mu_{\mathcal{S}}, a) &= \langle \mu_{\mathcal{S}}, R(\cdot, a) \rangle_{\mathcal{H}_{\mathcal{S}}} + \gamma \langle \mu_{O'|a; \mu_{\mathcal{S}}}, V^*(\mu_{S'}^{a, (\cdot)}) \rangle_{\mathcal{H}_{\mathcal{O}}}, \end{aligned}$$

4 Empirical Expression

Training samples are a set of $D_n = \{(\tilde{s}_i, \tilde{o}_i), \tilde{a}_i, \tilde{R}_i, (\tilde{s}'_i, \tilde{o}'_i)\}_{i=1}^n$ according to a POMDP. We assume that the true state samples $\{(\tilde{s}_i, \tilde{s}'_i)\}$ are available for training, but not during the test phase.

Belief Embedding Update Rule:

Empirical estimates $\hat{\mu}_{\mathcal{S}}$, $\hat{\mu}_{O'|a; \hat{\mu}_{\mathcal{S}}}$, $\hat{\mu}_{S'}^{a, o'}$ take respective forms $\hat{\mu}_{\mathcal{S}} = \Upsilon \alpha$, $\hat{\mu}_{O'|a; \hat{\mu}_{\mathcal{S}}} = \Phi \beta'_{a, \alpha}$, $\hat{\mu}_{S'}^{a, o'} = \Upsilon \alpha'_{a, o'}$. Update rule $\alpha \mapsto \beta'_{a, \alpha}$ is a linear transformation $\beta'_{a, \alpha} = L_{O|S, a} \alpha$ for all $a \in \mathcal{A}$ by $n \times n$ matrix

$$L_{O|S, a} = (G_{\mathcal{S}} + \varepsilon_S n I_n)^{-1} G_{S S'} (G_{(S, A)} + \varepsilon_{(S, A)} n I_n)^{-1} G_{(S, A)}(S, a) \quad (5)$$

with $G_{S S'} := \Upsilon^{\top} \Upsilon'$, $G_{(S, A)}(S, a) := D(\mathbf{k}_{\mathcal{A}}(a)) G_{\mathcal{S}}$, and $\mathbf{k}_{\mathcal{A}}(a) = \Psi^{\top} \psi(a)$. Update rule $\beta'_{a, \alpha} \mapsto \alpha'_{a, o'}$ is a transformation $\alpha'_{a, o'} = R_{S|O}(\beta'_{a, \alpha}) \mathbf{k}_{O}(o')$ with a non-negative vector $\beta'_{a, \alpha}$ and $n \times n$ matrix

$$R_{S|O}(\beta'_{a, \alpha}) = (D(\beta'_{a, \alpha}) G_{O} + \varepsilon_n I_n)^{-1} D(\beta'_{a, \alpha}). \quad (6)$$

Claim 2. Given samples D_n , the empirical expression of the kernel Bellman optimality equation (Claim 1) is

$$\hat{V}^*(\alpha) = \max_{a \in \mathcal{A}} \hat{Q}^*(\alpha, a), \quad \hat{\pi}^*(\alpha) = \arg \max_{a \in \mathcal{A}} \hat{Q}^*(\alpha, a), \quad \hat{Q}^*(\alpha, a) = \alpha^{\top} \mathbf{R}_a + \gamma \beta'_{a, \alpha}^{\top} \mathbf{V}^*(\alpha'_{a, \mathcal{O}_0}),$$

where $\mathbf{R}_a = (R(\tilde{s}_1, a), \dots, R(\tilde{s}_n, a))^{\top} \in \mathbb{R}^n$ is the reward vector on samples \mathcal{S}_0 for action a and $\hat{\mathbf{V}}^*(\alpha'_{a, \mathcal{O}_0}) = (\hat{V}^*(\alpha'_{a, \tilde{o}_1}), \dots, \hat{V}^*(\alpha'_{a, \tilde{o}_n}))^{\top} \in \mathbb{R}^n$ is the posterior belief value vector on samples \mathcal{O}_0 given action a .

Kernel Value Iteration Algorithm:

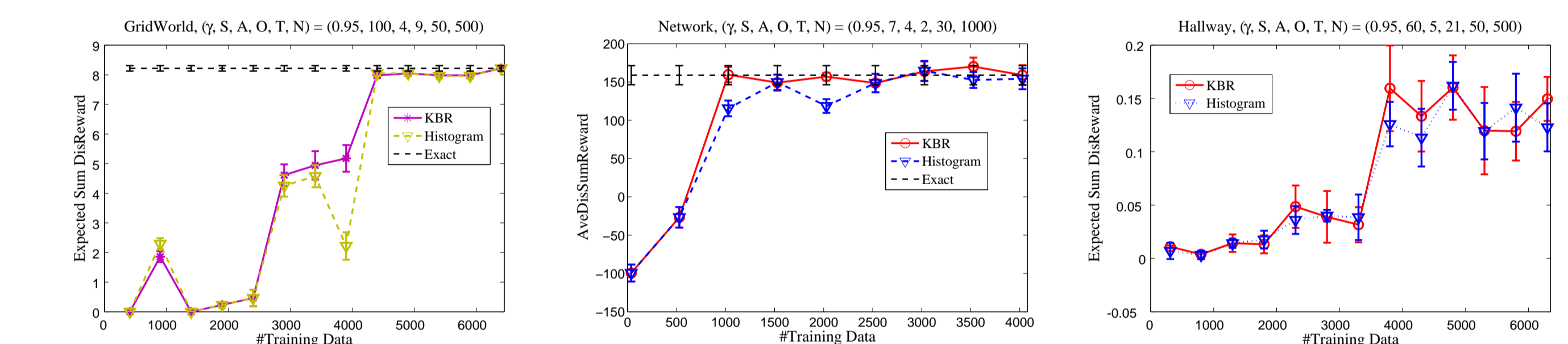
$\hat{V}_d = \hat{H}_n \hat{V}_{d-1}$ ($d \geq 1$), where \hat{H}_n is the kernel Bellman operator

$$(\hat{H}_n V)(\alpha) = \max_{a \in \mathcal{A}} \left[\alpha^{\top} \mathbf{R}_a + \gamma \beta'_{a, \alpha}^{\top} \mathbf{V}(\alpha'_{a, \mathcal{O}_0}) \right]. \quad (7)$$

\hat{H}_n can be enforced to be isotonic and contractive by replacing above weight vectors with probability vectors $\hat{\alpha}, \hat{\beta}'_{a, \alpha}, \hat{\alpha}'_{a, o'}$ as $\hat{w}_i = \frac{\max\{w_i, 0\}}{\sum_{i=1}^n \max\{w_i, 0\}}$ for weights w , proposed in [?].

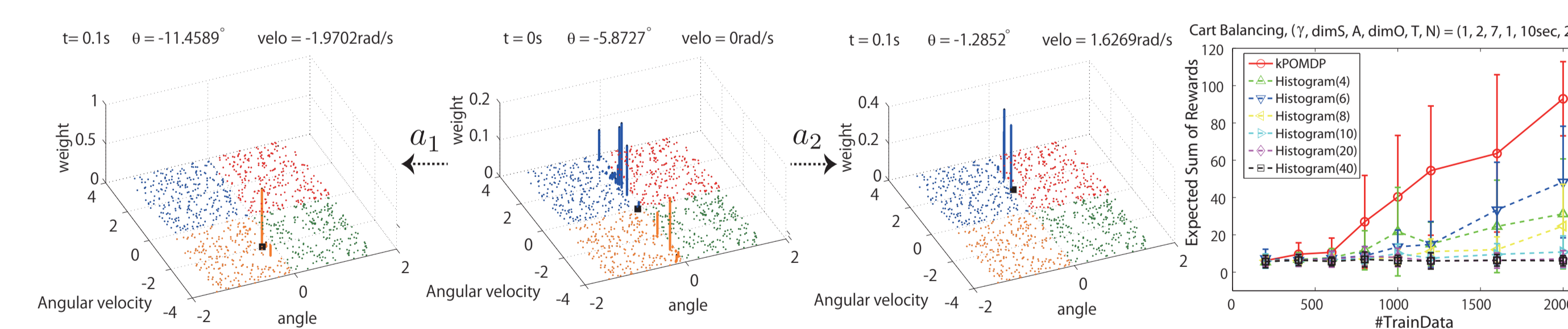
5 Experiments

5.1 Sets $\mathcal{S}, \mathcal{O}, \mathcal{A}$ are Finite



10×10 Gridworld (naive extension of 4×4 Gridworld), Network, and Hallway1 benchmark problems from left to right. Plots are averaged discounted sum of rewards earned in test experiments (vertical) against the number of training samples n (horizontal). Training samples are collected by uniform random actions. Parameters are within the title $(\gamma, S, A, O, T, N) = (\gamma, |\mathcal{S}|, |\mathcal{A}|, |\mathcal{O}|, T, N)$, where T and N indicates that one episode consists of T steps and results are averaged over N trials. kPOMDP are compared with histogram methods, in which transition and observation matrices are naively estimated by histograms.

5.2 Sets \mathcal{S}, \mathcal{O} are Euclidian spaces \mathbb{R}^d



An inverted pendulum problem, where hidden state s is angular and angular velocity $(\theta, \dot{\theta})$, and angular θ is only observed. All the points (colored by RGBY) in the three 3D plots indicate training samples on hidden states \mathcal{S} , and z axis indicates belief embedding weights $\hat{\alpha}$ on their samples (after normalization). The true hidden state is marked by the black point in each 3D plot. Since $\dot{\theta}$ is uncertain at the initial point, positive weights spread in the direction of $\dot{\theta}$ axis in the middle 3D figure. The left and right 3D figures show that θ is well estimated by the kPOMDPs dynamics for both of executed actions a_1 and a_2 .