

An Online Spectral Learning Algorithm for Partially Observable Nonlinear Dynamical Systems

Byron Boots

Machine Learning Department
Carnegie Mellon University
Pittsburgh, Pennsylvania 15213

Geoffrey J. Gordon

Machine Learning Department
Carnegie Mellon University
Pittsburgh, Pennsylvania 15213

Abstract

Recently, a number of researchers have proposed *spectral* algorithms for learning models of dynamical systems—for example, Hidden Markov Models (HMMs), Partially Observable Markov Decision Processes (POMDPs), and Transformed Predictive State Representations (TPSRs). These algorithms are attractive since they are statistically consistent and not subject to local optima. However, they are *batch* methods: they need to store their entire training data set in memory at once and operate on it as a large matrix, and so they cannot scale to extremely large data sets (either many examples or many features per example). In turn, this restriction limits their ability to learn accurate models of complex systems. To overcome these limitations, we propose a new *online* spectral algorithm, which uses tricks such as incremental Singular Value Decomposition (SVD) and random projections to scale to much larger data sets and more complex systems than previous methods. We demonstrate the new method on an inertial measurement prediction task and a high-bandwidth video mapping task and we illustrate desirable behaviors such as “closing the loop,” where the latent state representation changes suddenly as the learner recognizes that it has returned to a previously known place.

Introduction

Many problems in machine learning and artificial intelligence involve discrete-time partially observable nonlinear dynamical systems. If the observations are discrete, then *Hidden Markov Models* (HMMs) (Rabiner 1989) or, in the controlled setting, *Partially Observable Markov Decision Processes* (POMDPs) (Sondik 1971) can be used to represent belief as a discrete distribution over latent states. *Predictive State Representations* (PSRs) (Littman, Sutton, and Singh 2002), *Transformed Predictive State Representations* (TPSRs) (Rosencrantz, Gordon, and Thrun 2004; Boots, Siddiqi, and Gordon 2010), and the closely related *Observable Operator Models* (OOMs) (Jaeger 2000) are generalizations of POMDPs that have attracted interest because they both have greater representational capacity than

POMDPs and yield representations that are *at least* as compact (Singh, James, and Rudary 2004). In contrast to the latent-variable representations of POMDPs, predictive representations like PSRs, TPSRs, and OOMs represent the state of a dynamical system by tracking occurrence probabilities of a set of future events (called *tests* or *characteristic events*) conditioned on past events (called *histories* or *indicative events*).

Recently, spectral algorithms have been increasingly used to learn models of partially observable nonlinear dynamical systems such as HMMs (Hsu, Kakade, and Zhang 2009; Siddiqi, Boots, and Gordon 2010) and TPSRs (Rosencrantz, Gordon, and Thrun 2004; Boots, Siddiqi, and Gordon 2010; Boots and Gordon 2010). Most of these algorithms are *statistically consistent*, unlike the popular expectation maximization (EM) algorithm, which is subject to local optima. Furthermore, spectral learning algorithms are easy to implement with a series of linear algebra operations. Despite these attractive features, spectral algorithms have so far had an important drawback: they are *batch* methods (needing to store their entire training data set in memory at once) instead of *online* ones (with space complexity independent of the number of training examples and time complexity linear in the number of training examples).

To remedy this drawback, we propose a fast, online spectral algorithm for TPSRs. TPSRs *subsume* HMMs, PSRs, and POMDPs (Singh, James, and Rudary 2004; Rosencrantz, Gordon, and Thrun 2004). In fact, previous spectral learning algorithms for several types of HMMs (Hsu, Kakade, and Zhang 2009; Siddiqi, Boots, and Gordon 2010; Song et al. 2010) are more accurately described as TPSR learning algorithms *applied to* HMMs. Therefore, our algorithm also improves on past algorithms for these other models. Our method leverages fast, low-rank modifications of the thin singular value decomposition (Brand 2006), and uses tricks such as random projections to scale to extremely large numbers of examples and features per example. Consequently, the new method can handle orders of magnitude larger data sets than previous methods, and can therefore scale to learn models of systems that are too complex for previous methods.

Experiments show that our online spectral learning algorithm does a good job recovering the parameters of a nonlinear dynamical system in two partially observable do-

mains. In our first experiment we empirically demonstrate that our online spectral learning algorithm is unbiased by recovering the parameters of a small but difficult synthetic Reduced-Rank HMM. In our second experiment we demonstrate the performance of the new method on a difficult, high-bandwidth video understanding task.

Predictive State Representations

We take a Predictive State Representation (PSR) (Littman, Sutton, and Singh 2002) view of non-linear dynamical systems. A PSR is a compact description of a dynamical system that represents state as a set of predictions of observable experiments or *tests*. Specifically, a test is an ordered sequence of action-observation pairs $q = [a_1^q, o_1^q, \dots, a_k^q, o_k^q]$ that can be executed and observed at a given time. If the observations produced by the dynamical system match those specified by the test, the test is said to have *succeeded*. The *prediction* for q , $\mathbb{P}[q^O \mid \text{do}(q^A)]$, is the probability of seeing observations $q^O = [o_1^q, \dots, o_k^q]$, given that we *intervene* (Pearl 2000) to take the actions $q^A = [a_1^q, \dots, a_k^q]$. The key idea behind a PSR is that, if we know the expected outcomes of all possible tests, then we know everything there is to know about state.

A *history* is an ordered sequence of action-observation pairs $h = [a_1^h, o_1^h, \dots, a_t^h, o_t^h]$ that has been executed and observed prior to a given time. We write $\mathcal{Q}(h)$ for the *prediction vector* of success probabilities for a set of tests $\mathcal{Q} = \{q_i\}$ given a history h .

Formally a PSR consists of the tuple $\langle \mathcal{A}, \mathcal{O}, \mathcal{Q}, \mathcal{F}, m_1 \rangle$. \mathcal{A} is the set of possible actions, and \mathcal{O} is the set of possible observations. \mathcal{Q} is a *core* set of tests, i.e., a set such that $\mathcal{Q}(h)$ is a sufficient statistic for predicting *all* tests: the prediction for any test τ is a function of $\mathcal{Q}(h)$. \mathcal{F} is the set of functions f_τ which embody these predictions: $\mathbb{P}[\tau^O \mid h, \text{do}(\tau^A)] = f_\tau(\mathcal{Q}(h))$. Finally, $m_1 = \mathcal{Q}(\epsilon)$ is the initial prediction vector.

We restrict ourselves to *linear* PSRs, in which all prediction functions are linear: $f_\tau(\mathcal{Q}(h)) = r_\tau^\top \mathcal{Q}(h)$ for some vector $r_\tau \in \mathbb{R}^{|\mathcal{Q}|}$. A core set \mathcal{Q} for a linear PSR is said to be *minimal* if the tests in \mathcal{Q} are linearly independent (Jaeger 2000; Singh, James, and Rudary 2004), i.e., no one test’s prediction is a linear function of the other tests’ predictions. Note that the restriction to linear prediction functions is only a restriction to linear relationships between conditional probabilities of tests; linear PSRs can still represent systems with *nonlinear* dynamics.

Since $\mathcal{Q}(h)$ is a sufficient statistic for all test predictions, it is a *state* for our PSR: i.e., we can remember just $\mathcal{Q}(h)$ instead of h itself. After action a and observation o , we can update $\mathcal{Q}(h)$ recursively: if we write M_{ao} for the matrix with rows $r_{ao\tau}^\top$ for $\tau \in \mathcal{Q}$, then we can use Bayes’ Rule to show:

$$\mathcal{Q}(hao) = \frac{M_{ao}\mathcal{Q}(h)}{\mathbb{P}[o \mid h, \text{do}(a)]} = \frac{M_{ao}\mathcal{Q}(h)}{m_\infty^\top M_{ao}\mathcal{Q}(h)} \quad (1)$$

where m_∞ is defined by $m_\infty^\top \mathcal{Q}(h) = 1 (\forall h)$.

Transformed PSRs

Transformed PSRs (TPSRs) (Rosencrantz, Gordon, and Thrun 2004; Boots, Siddiqi, and Gordon 2010) are a gen-

eralization of PSRs: for any invertible matrix J , if the parameters m_1 , M_{ao} , and m_∞ represent a PSR, then the transformed parameters $b_1 = Jm_1$, $B_{ao} = JM_{ao}J^{-1}$, and $b_\infty = J^{-\top}m_\infty$ represent an *equivalent* TPSR. In addition to the initial TPSR state b_1 , we define normalized conditional *internal states* b_t , which we can update similarly to Eq. 1:

$$b_{t+1} \equiv \frac{B_{ao_{1:t}}b_1}{b_\infty^\top B_{ao_{1:t}}b_1} = \frac{B_{a_t o_t}b_t}{b_\infty^\top B_{a_t o_t}b_t} \quad (2)$$

Pairs $J^{-1}J$ cancel during the update, showing that predictions are equivalent as claimed:

$$\begin{aligned} \Pr[o_{1:t} \mid \text{do}(a_{1:t})] &= m_\infty^\top M_{ao_{1:t}} m_1 \\ &= m_\infty^\top J^{-1} J M_{ao_{1:t}} J^{-1} m_1 \\ &= b_\infty^\top B_{ao_{1:t}} b_1 \end{aligned} \quad (3)$$

By choosing the invertible transform J appropriately (see the next subsection), we can think of TPSRs as maintaining a small number of sufficient statistics which are *linear combinations* of predictions for a (potentially very large) core set of tests. This view leads to the main benefit of TPSRs over regular PSRs: given a core set of tests, we can find low dimensional parameters using spectral methods and regression instead of combinatorial search. In this respect, TPSRs are closely related to the transformed representations of LDSs and HMMs found by *subspace identification* (Van Overschee and De Moor 1996; Katayama 2005; Soatto and Chiuso 2001; Hsu, Kakade, and Zhang 2009). Furthermore, to make it practical to work with data gathered from complex real-world systems, we can learn from finite-dimensional features of the past and future, rather than an extremely large or even infinite core set of tests (see Section “Batch Learning of TPSRs,” below). Additional details regarding the relationship between TPSRs and PSRs can be found in (Boots, Siddiqi, and Gordon 2010).

Relating TPSRs and PSRs For some PSR, let \mathcal{Q} be a minimal core set of tests. Then, let \mathcal{T} be a (larger) core set of tests, and let \mathcal{H} be a mutually exclusive and exhaustive partition of the set of all possible histories. (Elements of \mathcal{H} are called *indicative events* (Jaeger 2000).) And, let \mathcal{AO} be the set of all possible action-observation pairs. Define $\phi^{\mathcal{H}}(h)$ for $h \in \mathcal{H}$ to be a vector of *indicative* features, i.e., features of history, and define $\phi^{AO}(a, o)$ to be a vector of features of a present action and observation. Finally, define $\phi^{\mathcal{T}}(h)$ to be a vector of *characteristic* features: that is, each entry of $\phi^{\mathcal{T}}(h)$ is a linear combination of some set of test predictions.

For purposes of gathering data, we assume that we can sample from a sufficiently diverse distribution ω over histories. Note that this assumption means that we cannot estimate m_1 (equivalently, b_1), since we typically don’t have samples of trajectories starting from m_1 . Instead, we will estimate m_* , an arbitrary feasible state. If we only have m_* , initial probability estimates will be approximate, but the difference will disappear over time as our process mixes. We will also assume that we execute a known exploration policy from each sampled history; with this assumption, it is possible to construct unbiased samples of $\phi^{\mathcal{T}}(h)$ by importance weighting (Bowling et al. 2006; Boots and Gordon 2010).

When our algorithms below call for samples of $\phi^T(h)$, we use this importance weighting trick to provide them.

We define Φ^T , Φ^H , and Φ^{AO} as matrices of characteristic, indicative, and present features respectively, with first dimension equal to the number of features and second dimension equal to $|\mathcal{H}|$. An entry of Φ^H is the expectation of one of the indicative features given the occurrence of one of the indicative events and the exploration policy; an entry of Φ^T is the expectation of one of our characteristic features given one of the indicative events; and an entry of Φ^{AO} is the expectation of one of the present features given one of the indicative events and the exploration policy. We also define $\psi = \mathbb{P}[\mathcal{H}]$, $D = \text{diag}(\psi)$, $R \in \mathbb{R}^{|\mathcal{T}| \times |\mathcal{Q}|}$ as the matrix with rows $r_{\tau_i}^T$, $S \in \mathbb{R}^{|\mathcal{Q}| \times |\mathcal{H}|}$ as the expected state $\mathbb{E}[\mathcal{Q} | \mathcal{H}]$, and M as a $|\mathcal{Q}| \times |\mathcal{AO}| \times |\mathcal{Q}|$ third-order tensor (each $|\mathcal{Q}| \times |\mathcal{Q}|$ slice, M_{ao} , of M is the transition matrix for a particular action-observation pair).

Given the above notation, we define several covariance and “trivariance” matrices which are related to the parameters of the PSR. In several of the following equations we use *tensor-matrix multiplication* \times_v , also known as a *mode- v product*: \times_v multiplies the second dimension of a matrix with the v th mode of a tensor.

$$\begin{aligned} [\mu_{\mathcal{H}}]_i &\equiv \mathbb{E}[\phi_i^H(h)] \\ &\implies \mu_{\mathcal{H}} = \Phi^H \psi \end{aligned} \quad (4a)$$

$$\begin{aligned} [\Sigma_{\mathcal{AO}, \mathcal{AO}}]_{i,j} &\equiv \mathbb{E}[\phi_i^{AO}(a, o) \cdot \phi_j^{AO}(a, o)] \\ \implies \Sigma_{\mathcal{AO}, \mathcal{AO}} &= \Phi^{AO} D \Phi^{AO T} \end{aligned} \quad (4b)$$

$$\begin{aligned} [\Sigma_{\mathcal{T}, \mathcal{H}}]_{i,j} &\equiv \mathbb{E}[\phi_i^T(\tau^O) \cdot \phi_j^H(h) | \text{do}(\tau^A)] \\ \implies \Sigma_{\mathcal{T}, \mathcal{H}} &= \Phi^T R S D \Phi^H T \end{aligned} \quad (4c)$$

$$\begin{aligned} [\Sigma_{\mathcal{T}, \mathcal{AO}, \mathcal{H}}]_{i,j,k} &\equiv \mathbb{E}[\phi_i^T(\tau^O) \cdot \phi_j^H(h) \cdot \phi_k^{AO}(ao) | \text{do}(\tau^A, a)] \\ \implies \Sigma_{\mathcal{T}, \mathcal{AO}, \mathcal{H}} &= M \times_1 (\Phi^T R) \times_2 (\Phi^{AO} D) \times_3 (\Phi^H D S^T) \end{aligned} \quad (4d)$$

Now, if we are given an additional matrix U such that $U^T \Phi^T R$ is invertible, we can use Equations 4a–d to define a TPSR whose parameters are only a similarity transform away from the original PSR parameters.

$$b_* \equiv U^T \Sigma_{\mathcal{T}, \mathcal{H}} e = (U^T \Phi^T R) m_* \quad (5a)$$

$$b_{\infty}^T \equiv \mu_{\mathcal{H}}^T (U^T \Sigma_{\mathcal{T}, \mathcal{H}})^{\dagger} = m_{\infty}^T (U^T \Phi^T R)^{-1} \quad (5b)$$

$$\begin{aligned} B_{ao} &\equiv \Sigma_{\mathcal{T}, \mathcal{AO}, \mathcal{H}} \\ &\times_1 U^T \times_2 \Phi^{AO T} (\Sigma_{\mathcal{AO}, \mathcal{AO}})^{-1} \times_3 (\Sigma_{\mathcal{T}, \mathcal{H}} U)^{\dagger} \\ &= (U^T \Phi^T R) M_{ao} (U^T \Phi^T R)^{-1} \end{aligned} \quad (5c)$$

Batch Learning of TPSRs

The identities in Equation 5a–c yield a straightforward learning algorithm (Boots, Siddiqi, and Gordon 2010): we build empirical estimates $\hat{\mu}_{\mathcal{H}}$, $\hat{\Sigma}_{\mathcal{AO}, \mathcal{AO}}$, $\hat{\Sigma}_{\mathcal{T}, \mathcal{H}}$, and $\hat{\Sigma}_{\mathcal{T}, \mathcal{AO}, \mathcal{H}}$ of the matrices defined in Equation 4. Once we have constructed $\hat{\Sigma}_{\mathcal{T}, \mathcal{H}}$, we can compute \hat{U} as the matrix of n leading left singular vectors of $\hat{\Sigma}_{\mathcal{T}, \mathcal{H}}$. Finally, we plug the estimated covariances and \hat{U} into Equation 5 to compute

estimated PSR parameters. One of the advantages of subspace identification is that the complexity of the model can be tuned by selecting the number of singular vectors in \hat{U} , at the risk of losing prediction quality.

As we include more data in our averages, the law of large numbers guarantees that our estimates converge to the true matrices $\mu_{\mathcal{H}}$, $\Sigma_{\mathcal{AO}, \mathcal{AO}}$, $\Sigma_{\mathcal{T}, \mathcal{H}}$, and $\Sigma_{\mathcal{T}, \mathcal{AO}, \mathcal{H}}$. So by continuity of the formulas above, if our system is truly a PSR of finite rank, our estimates \hat{b}_* , \hat{b}_{∞} , and \hat{B}_{ao} converge, with probability 1, to the true parameters up to a linear transform—that is, our learning algorithm is *consistent*.¹

An Efficient Batch Learning Algorithm

Unfortunately, it is difficult to implement the naïve algorithm in practice. For example, if there are a large number of features of tests or features of histories, we may not be able even to store the the tensor in Equation 4d. Therefore, we need to use a more efficient algorithm, one that does not explicitly build the tensor $\Sigma_{\mathcal{T}, \mathcal{AO}, \mathcal{H}}$.

The key idea is to compute a set of smaller-sized intermediate quantities from realizations of characteristic features ϕ^T , indicative features ϕ^H , and present features ϕ^{AO} , and then compute TPSR parameters from these quantities. In particular, we compute a subset of the empirical estimates from above: $\hat{\mu}_{\mathcal{H}}$, $\hat{\Sigma}_{\mathcal{AO}, \mathcal{AO}}$ and $\hat{\Sigma}_{\mathcal{T}, \mathcal{H}}$. From $\hat{\Sigma}_{\mathcal{T}, \mathcal{H}}$ we compute the rank- n singular value decomposition \hat{U} , \hat{S} , \hat{V}^T . Then, instead of computing $\hat{\Sigma}_{\mathcal{T}, \mathcal{AO}, \mathcal{H}}$, we use the above matrices to compute a smaller tensor $\hat{B}_{\mathcal{AO}} = \hat{\Sigma}_{\mathcal{T}, \mathcal{AO}, \mathcal{H}} \times_1 \hat{U}^T \times_3 (\hat{\Sigma}_{\mathcal{T}, \mathcal{H}})^{\dagger}$ directly. To get the final desired parameter \hat{B}_{ao} (Equation 5), we simply compute $B_{ao} = \hat{B}_{\mathcal{AO}} \times_2 \phi^{AO}(ao)^T (\Sigma_{\mathcal{AO}, \mathcal{AO}})^{-1}$. The tensor $\hat{B}_{\mathcal{AO}}$ has much lower dimensionality than $\Sigma_{\mathcal{T}, \mathcal{AO}, \mathcal{H}}$, with first and third modes being no greater than n dimensions (the length of the latent state vector b_t). Therefore, we save substantially on both memory and runtime by avoiding construction of the larger tensor.

In more detail, we begin by estimating $\mu_{\mathcal{H}}$, the unnormalized² empirical expectation (over histories sampled from ω) of the indicative feature vector: if ϕ_t^H is the feature vector for the t th history,

$$\hat{\mu}_{\mathcal{H}} = \sum_{t=1}^w \phi_t^H \quad (6a)$$

Next, we estimate $\Sigma_{\mathcal{AO}, \mathcal{AO}}^{-1}$, the inverse of the unnormalized empirical expectation of the product of present features:

$$\hat{\Sigma}_{\mathcal{AO}, \mathcal{AO}}^{-1} = \left(\sum_{t=1}^w \phi_t^{AO} \phi_t^{AO} \right)^{-1} \quad (6b)$$

Next, each element of our estimate $\hat{\Sigma}_{\mathcal{T}, \mathcal{H}}$ is an unnormalized empirical expectation of the *product* of one indicative

¹Continuity holds if we fix \hat{U} ; a similar but more involved argument works if we estimate \hat{U} as well.

²We can get away with using *unnormalized* empirical expectations in Equation 6a–d since the normalizations would just cancel when we compute the TPSR parameters in Equation 7a–c.

feature and one characteristic feature, if we sample a history from ω and then follow an appropriate sequence of actions. We can compute all elements of $\widehat{\Sigma}_{\mathcal{T},\mathcal{H}}$ from a single sample of trajectories if we sample histories from ω , follow an appropriate exploratory policy, and then importance-weight each sample (Bowling et al. 2006): $[\widehat{\Sigma}_{\mathcal{T},\mathcal{H}}]_{ij}$ is $\sum_{t=1}^w \lambda_t \phi_{it}^{\mathcal{T}} \phi_{jt}^{\mathcal{H}}$, where λ_t is an importance weight.

Next we compute $\widehat{U}\widehat{S}\widehat{V}^T$, the rank- n thin singular value decomposition of $\widehat{\Sigma}_{\mathcal{T},\mathcal{H}}$:

$$\langle \widehat{U}, \widehat{S}, \widehat{V} \rangle = \text{SVD} \left(\sum_{t=1}^w \lambda_t \phi_t^{\mathcal{T}} \phi_t^{\mathcal{H}} \right) \quad (6c)$$

The left singular vectors \widehat{U} are directly needed, e.g., in Eq. 5. The right singular vectors \widehat{V} and singular values \widehat{S} can also be used to make computation of the other TPSR parameters more efficient: recall from Equation 4d that each element of $\widehat{\Sigma}_{\mathcal{T},ao,\mathcal{H}}$ is a weighted unnormalized empirical expectation of the product of one indicative feature, one characteristic feature, and one present feature. To efficiently construct the tensor \widehat{B}_{AO} , we combine Equation 4d with Equation 5c while making use of the singular value decomposition of $\widehat{\Sigma}_{\mathcal{T},\mathcal{H}}$:

$$\widehat{B}_{AO} = \sum_{t=1}^w \lambda_t \left(\widehat{U}^T \phi_t^{\mathcal{T}} \right) \otimes \left(\phi_t^{AO} \right) \otimes \left(\widehat{S}^{-1} \widehat{V}^T \phi_t^{\mathcal{H}} \right) \quad (6d)$$

where \otimes indicates tensor (outer) product. The idea here is to build the lower-dimensional \widehat{B}_{AO} directly, by making use of the tall thin matrices \widehat{U} and \widehat{V} , without first building the high-dimensional $\widehat{\Sigma}_{\mathcal{T},AO,\mathcal{H}}$.

Using these matrices we can efficiently compute estimates of the TPSR parameters in Equation 5:

$$\hat{b}_* = \widehat{S}\widehat{V}^T e \quad (7a)$$

$$\hat{b}_\infty^T = \widehat{\mu}_{\mathcal{H}}^T \widehat{V} \widehat{S}^{-1} \quad (7b)$$

$$\widehat{B}_{ao} = \widehat{B}_{AO} \times_2 \left(\phi^{AO}(ao)^T \cdot \widehat{\Sigma}_{AO,AO}^{-1} \right) \quad (7c)$$

This learning algorithm works well when the number of features of tests, histories, and action-observation pairs is relatively small, and in cases where data is collected in batch. These restrictions can be limiting for many real-world data sets. In practice, the number of features may need to be quite large in order to accurately estimate the parameters of the TPSR. Additionally, we are often interested in estimating TPSRs from massive datasets, updating TPSR parameters given a new batch of data, or learning TPSRs online from a data stream. Below we develop several computationally efficient extensions to overcome these practical obstacles to learning in real-world situations.

Iterative Updates to TPSR Parameters

We first attack the problem of updating existing TPSR parameters given a batch of new information. Next, we look at the special case of updating TPSR parameters in an online setting (batch size 1), and develop additional optimizations for this situation.

Batch Updates We first present an algorithm for updating existing TPSR parameters given a new batch of characteristic features $\phi_{\text{new}}^{\mathcal{T}}$, indicative features $\phi_{\text{new}}^{\mathcal{H}}$, and present features ϕ_{new}^{AO} . Naïvely, we can just store empirical estimates and update them from each new batch of data: $\hat{\mu}_{\mathcal{H}} + \sum_{t=1}^w \phi_{\text{new},t}^{\mathcal{H}}$, $\widehat{\Sigma}_{AO,AO} + \phi_{\text{new}}^{AO} \phi_{\text{new}}^{AO^T}$, $\widehat{\Sigma}_{\mathcal{T},\mathcal{H}} + \phi_{\text{new}}^{\mathcal{T}} \phi_{\text{new}}^{\mathcal{H}^T}$, and $\widehat{\Sigma}_{\mathcal{T},AO,\mathcal{H}} + \sum_{t=1}^w \phi_{\text{new},t}^{\mathcal{T}} \otimes \phi_{\text{new},t}^{AO} \otimes \phi_{\text{new},t}^{\mathcal{H}}$. Then, after each batch, we can apply Equation 5a–c to learn new TPSR parameters \hat{b}_* , \hat{b}_∞^T , and \widehat{B}_{ao} .

This naïve algorithm is very inefficient: it requires storing each of the matrices in Equation 4a–d, updating these matrices given new information, and recomputing the TPSR parameters. However, as we have seen, it is also possible to write the TPSR parameters (Equation 7a–c) in terms of a set of lower-dimensional memory-efficient matrices and tensors (Equation 6a–d), made possible by the singular value decomposition of $\widehat{\Sigma}_{\mathcal{T},\mathcal{H}}$. The key idea is to update these lower-dimensional matrices directly, instead of the naïve updates suggested above, by taking advantage of numerical algorithms for updating singular value decompositions efficiently (Brand 2006).

We begin by simply implementing the additive update to the vector $\hat{\mu}_{\mathcal{H}}$.

$$\hat{\mu}_{\mathcal{H}_{\text{new}}} = \hat{\mu}_{\mathcal{H}} + \sum_{t=1}^w \phi_t^{\mathcal{H}} \quad (8a)$$

The inverse empirical covariance of present features remains computationally expensive to update. If the new batch of data is large, and updates infrequent, then we can maintain the empirical covariance $\widehat{\Sigma}_{AO,AO}$ separately from the inverse. We update $\widehat{\Sigma}_{AO,AO}$, and after each update, we invert the matrix.

$$\widehat{\Sigma}_{AO,AO_{\text{new}}}^{-1} = \left(\widehat{\Sigma}_{AO,AO} + \phi_{\text{new}}^{AO} \phi_{\text{new}}^{AO^T} \right)^{-1} \quad (8b)$$

If we are updating frequently with smaller batches of data, we can instead use a low-rank update via the Sherman-Morrison formula. See Equation 9a for details.

The main computational savings come from using incremental SVD to update \widehat{U} , \widehat{S} , \widehat{V} and \widehat{B}_{AO} directly. The incremental update for \widehat{U} , \widehat{S} , \widehat{V} is much more efficient than the naïve additive update when the number of new data points is much smaller than the number of features in $\phi^{\mathcal{T}}$ and $\phi^{\mathcal{H}}$. The incremental update for \widehat{B}_{AO} saves time and space when the number of features in $\phi^{\mathcal{T}}$ and $\phi^{\mathcal{H}}$ is much larger than the latent dimension n .

Our goal is therefore to compute the updated SVD,

$$\langle \widehat{U}_{\text{new}}, \widehat{S}_{\text{new}}, \widehat{V}_{\text{new}} \rangle = \text{SVD} \left(\widehat{\Sigma}_{\mathcal{T},\mathcal{H}} + \phi_{\text{new}}^{\mathcal{T}} \Lambda \phi_{\text{new}}^{\mathcal{H}^T} \right),$$

where Λ is a diagonal matrix of importance weights $\Lambda = \text{diag}(\lambda_{1:t})$. We will derive the incremental SVD update in two steps. First, if the new data $\phi_{\text{new}}^{\mathcal{T}}$ and $\phi_{\text{new}}^{\mathcal{H}}$ lies entirely within the column spaces of \widehat{U} and \widehat{V} respectively, then we can find \widehat{S}_{new} by projecting both the new and old data onto the subspaces defined by \widehat{U} and \widehat{V} , and diagonalizing the

resulting small ($n \times n$) matrix:

$$\begin{aligned} \langle \hat{U}, \hat{S}_{\text{new}}, \hat{V} \rangle &= \text{SVD} \left(\hat{U}^T \left(\hat{\Sigma}_{\mathcal{T}, \mathcal{H}} + \phi_{\text{new}}^T \Lambda \phi_{\text{new}}^{\mathcal{H}T} \right) \hat{V} \right) \\ &= \text{SVD} \left(\hat{S} + (\hat{U}^T \phi_{\text{new}}^T) \Lambda (\hat{V}^T \phi_{\text{new}}^{\mathcal{H}T}) \right) \end{aligned}$$

We can then compute $\hat{U}_{\text{new}} = \hat{U} \hat{U}^T$ and $\hat{V}_{\text{new}} = \hat{V} \hat{V}^T$, the rotations of \hat{U} and \hat{V} induced by the new data.

If the new data does *not* lie entirely within the column space of \hat{U} and \hat{V} , we can update the SVD efficiently (and optionally approximately) following Brand (2006). The idea is to split the new data into a part within the column span of \hat{U} and \hat{V} and a remainder, and use this decomposition to construct a small matrix to diagonalize as above.

Let C and D be orthonormal bases for the component of the column space of ϕ_{new}^T orthogonal to \hat{U} and the component of the column space of $\phi_{\text{new}}^{\mathcal{H}}$ orthogonal to \hat{V} :

$$\begin{aligned} C &= \text{orth} \left((I - \hat{U} \hat{U}^T) \phi_{\text{new}}^T \right) \\ D &= \text{orth} \left((I - \hat{V} \hat{V}^T) \phi_{\text{new}}^{\mathcal{H}} \right) \end{aligned}$$

The dimension of C and D is upper-bounded by the number of data points in our new batch, or the number of features of tests and histories, whichever is smaller. (If the dimension is large, the orthogonalization step above (as well as other steps below) may be too expensive; we can accommodate this case by splitting a large batch of examples into several smaller batches.) Let

$$K = \begin{bmatrix} \hat{S} & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} \hat{U}^T \\ C^T \end{bmatrix} \phi_{\text{new}}^T \Lambda \phi_{\text{new}}^{\mathcal{H}T} \begin{bmatrix} \hat{V} \\ D \end{bmatrix},$$

and diagonalize K to get the update to \hat{S} :

$$\langle \hat{U}, \hat{S}_{\text{new}}, \hat{V} \rangle = \text{SVD}(K) \quad (8c)$$

Finally, as above, \hat{U} and \hat{V} rotate the extended subspaces $[\hat{U} \ C]$ and $[\hat{V} \ D]$:

$$\hat{U}_{\text{new}} = [\hat{U} \ C] \hat{U} \quad (8d)$$

$$\hat{V}_{\text{new}} = [\hat{V} \ D] \hat{V} \quad (8e)$$

Note that if there are components orthogonal to \hat{U} and \hat{V} in the new data (i.e., if C and D are nonempty), the size of the thin SVD will grow. So, during this step, we may choose to tune the complexity of our estimated model by restricting the dimensionality of the SVD. If we do so, we may *lose information* compared to a batch SVD: if future data causes our estimated leading singular vectors to change, the dropped singular vectors may become relevant again. However, empirically, this information loss can be minimal, especially if we keep extra ‘‘buffer’’ singular vectors beyond what we expect to need.

Also note that the above updates do not necessarily preserve orthonormality of \hat{U}_{new} and \hat{V}_{new} , due to discarded nonzero singular values and the accumulation of numerical

errors. To correct for this, every few hundred iterations, we re-orthonormalize using a QR decomposition and a SVD:

$$\langle U_Q, U_R \rangle = \text{QR} \left(\hat{U}_{\text{new}} \right)$$

$$\langle V_Q, V_R \rangle = \text{QR} \left(\hat{V}_{\text{new}} \right)$$

$$\langle U_{QR}, S_{QR}, V_{QR} \rangle = \text{SVD} \left(U_R \hat{S}_{\text{new}} V_R^T \right)$$

$$\hat{U}_{\text{new}} = U_Q U_{QR}$$

$$\hat{V}_{\text{new}} = V_Q V_{QR}$$

$$\hat{S}_{\text{new}} = S_{QR}$$

The updated SVD now gives us enough information to compute the update to $\hat{B}_{\mathcal{A}\mathcal{O}}$. Let Λ be the diagonal tensor of importance weights $\Lambda_{i,i,i} = \lambda_i (i \in 1, 2, \dots, t)$. Using the newly computed subspaces \hat{U}_{new} , \hat{S}_{new} , and \hat{V}_{new} , we can compute an additive update from the new data.

$$\begin{aligned} \hat{B}_{\mathcal{A}\mathcal{O}\text{new}} &= \left(\hat{\Sigma}_{\mathcal{T}, \mathcal{A}\mathcal{O}, \mathcal{H}} + \Lambda \times_1 \phi_{\text{new}}^T \times_2 \phi_{\text{new}}^{\mathcal{A}\mathcal{O}} \times_3 \phi_{\text{new}}^{\mathcal{H}} \right) \\ &\quad \times_1 \hat{U}_{\text{new}}^T \times_3 \hat{S}_{\text{new}}^{-T} \hat{V}_{\text{new}}^T \\ &= \hat{B}_{\text{update}} + \hat{B}_{\text{new}} \end{aligned} \quad (8f)$$

where \hat{B}_{update} can be viewed as the projection of $\hat{B}_{\mathcal{A}\mathcal{O}}$ onto the new subspace:

$$\begin{aligned} \hat{B}_{\text{update}} &= \hat{\Sigma}_{\mathcal{T}, \mathcal{A}\mathcal{O}, \mathcal{H}} \times_1 \hat{U}_{\text{new}}^T \times_3 \hat{S}_{\text{new}}^{-T} \hat{V}_{\text{new}}^T \\ &= \begin{bmatrix} \hat{B}_{\mathcal{A}\mathcal{O}} & 0 \\ 0 & 0 \end{bmatrix} \times_1 \left(\hat{U}_{\text{new}}^T \begin{bmatrix} \hat{U} & 0 \end{bmatrix} \right) \\ &\quad \times_3 \left(\hat{S}_{\text{new}}^{-1} \hat{V}_{\text{new}}^T \begin{bmatrix} \hat{V} \hat{S} & 0 \end{bmatrix} \right) \end{aligned}$$

and \hat{B}_{new} is the projection of the additive update onto the new subspace:

$$\begin{aligned} \hat{B}_{\text{new}} &= \Lambda \times_1 (\hat{U}_{\text{new}}^T \phi_{\text{new}}^T) \times_2 (\phi_{\text{new}}^{\mathcal{A}\mathcal{O}}) \times_3 (\hat{S}_{\text{new}}^{-1} \hat{V}_{\text{new}}^T \phi_{\text{new}}^{\mathcal{H}}) \end{aligned}$$

Once the updated estimates in Equation 8a–f have been calculated, we can compute the new TPSR parameters by Equation 7a–c.

Online updates In the online setting (with just one sample per batch), the updates to $\hat{\Sigma}_{\mathcal{A}\mathcal{O}, \mathcal{A}\mathcal{O}}$ and \hat{S} are rank-1, allowing some additional efficiencies. The inverse empirical covariance $\hat{\Sigma}_{\mathcal{A}\mathcal{O}, \mathcal{A}\mathcal{O}}^{-1}$ can be updated directly using the Sherman-Morrison formula:

$$\hat{\Sigma}_{\mathcal{A}\mathcal{O}, \mathcal{A}\mathcal{O}\text{new}}^{-1} = \hat{\Sigma}_{\mathcal{A}\mathcal{O}, \mathcal{A}\mathcal{O}}^{-1} - \frac{\hat{\Sigma}_{\mathcal{A}\mathcal{O}, \mathcal{A}\mathcal{O}}^{-1} \phi_1^{\mathcal{A}\mathcal{O}} \phi_1^{\mathcal{A}\mathcal{O}T} \hat{\Sigma}_{\mathcal{A}\mathcal{O}, \mathcal{A}\mathcal{O}}^{-1}}{1 + \phi_1^{\mathcal{A}\mathcal{O}T} \hat{\Sigma}_{\mathcal{A}\mathcal{O}, \mathcal{A}\mathcal{O}}^{-1} \phi_1^{\mathcal{A}\mathcal{O}}} \quad (9a)$$

Next we can compute the rank-1 update to the matrices \hat{U} , \hat{S} , and \hat{V}^T . We can compute C and D efficiently via a simplified Gram-Schmidt step (Brand 2006):

$$\tilde{C} = (I - \hat{U} \hat{U}^T) \phi_1^T$$

$$C = \tilde{C} / \|\tilde{C}\|$$

$$\tilde{D} = (I - \hat{V} \hat{V}^T) \phi_1^{\mathcal{H}}$$

$$D = \tilde{D} / \|\tilde{D}\|$$

Finally, we can compute K by adding a rank-1 matrix to \hat{S} :

$$K = \begin{bmatrix} \hat{S} & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} \hat{U}^T \\ C^T \end{bmatrix} \phi_1^T \lambda_1 \phi_1^T [\hat{V} D],$$

We then compute the updated parameters using Eqs. 8c–e.

Random Projections for High Dimensional Feature Spaces

Despite their simplicity and wide applicability, HMMs, POMDPs, and PSRs are limited in that they are usually restricted to discrete observations, and the state is usually restricted to have only moderate cardinality. In Section “Transformed PSRs,” above, we described a feature-based representation for TPSRs that relaxes this restriction. Recently, Song et al. went a step further, and proposed a spectral learning algorithm for HMMs with continuous observations by representing distributions over these observations and continuous latent states as embeddings in an infinite dimensional Hilbert space (Song et al. 2010). These Hilbert Space Embeddings of HMMs (HSE-HMMs) use essentially the same framework as other spectral learning algorithms for HMMs and PSRs, but avoid working in the infinite-dimensional Hilbert space by the well-known “kernel trick.”

HSE-HMMs have been shown to perform well on several real-world datasets, often beating the next best method by a substantial margin. However, they scale poorly due to the need to work with the kernel matrix, whose size is quadratic in the number of training points.

We can overcome this scaling problem and learn TPSRs that approximate HSE-HMMs using *random features* for kernel machines (Rahimi and Recht 2007): we construct a large but finite set of random features which let us *approximate* a desired kernel using ordinary dot products. Rahimi and Recht show how to approximate several popular kernels, including radial basis function (RBF) kernels and Laplacian kernels.) The benefit of random features is that we can use fast linear methods that do not depend on the number of data points to approximate the original kernel machine.

HSE-HMMs are no exception: using random features of tests and histories, we can *approximate* a HSE-HMM with a TPSR. If we combine random features with the above online learning algorithm, we can approximate an HSE-HMM very closely by using an extremely large number of random features. Such a large set of features would overwhelm batch spectral learning algorithms, but our online method allows us to approximate an HSE-HMM very closely, and scale HSE-HMMs to orders of magnitude larger training sets or even to *streaming* datasets with an inexhaustible supply of training data.

Experimental Results

We designed 3 sets of experiments to evaluate the statistical properties and practical potential of our online spectral learning algorithm. In the first experiment we show the convergence behavior of the algorithm. In the second experiment we show how online spectral learning combined with random projections can be used to learn a TPSR that closely approximates the performance of a HSE-HMM. In the third

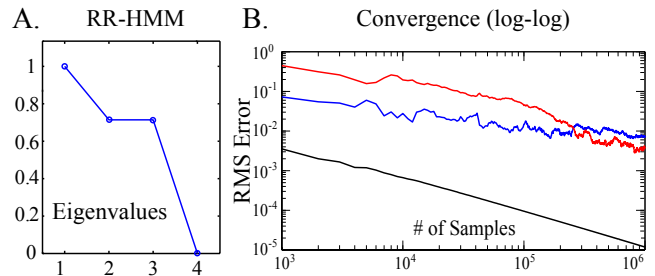


Figure 1: A synthetic RR-HMM. (A.) The eigenvalues of the true transition matrix. (B.) RMS error in the nonzero eigenvalues of the estimated transition matrix vs. number of training samples, averaged over 10 trials. The error steadily decreases, indicating that the TPSR model is becoming more accurate, as we incorporate more training data.

experiment we demonstrate how this combination allows us to model a high-bandwidth, high-dimensional video, where the amount of training data would overwhelm a kernel-based method like HSE-HMMs and the number of features would overwhelm a TPSR batch learning algorithm.

A Synthetic Example

First we demonstrate the convergence behavior of our algorithm on a difficult synthetic HMM from Siddiqi et al. (2010). This HMM is 2-step observable, with 4 states, 2 observations, and a rank-3 transition matrix. (So, the HMM is reduced rank (an “RR-HMM”) and features of multiple observations are required to disambiguate state.) The transition matrix T and the observation matrix O are:

$$T = \begin{bmatrix} 0.7829 & 0.1036 & 0.0399 & 0.0736 \\ 0.1036 & 0.4237 & 0.4262 & 0.0465 \\ 0.0399 & 0.4262 & 0.4380 & 0.0959 \\ 0.0736 & 0.0465 & 0.0959 & 0.7840 \end{bmatrix}$$

$$O = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

We sample observations from the true model and then estimate the model using the algorithm of Section “Online updates.” Since we only expect to recover the transition matrix up to a similarity transform, we compare the eigenvalues of $\hat{B} = \sum_o \hat{B}_o$ in the learned model to the eigenvalues of the transition matrix T of the true model. Fig. 1 shows that the learned eigenvalues converge to the true ones as the amount of data increases.

Slot Car Inertial Measurement

In a second experiment, we compare the online spectral algorithm with random features to HSE-HMMs with Gaussian RBF kernels. The setup consisted of a track and a miniature car (1:32 scale) guided by a slot cut into the track (Song et al. 2010). Figure 2(A) shows the car and the attached IMU (an Intel Inertiadot), as well as the 14m track, which contains elevation changes and banked curves. We collected the estimated 3D acceleration and velocity of the car at 10Hz. The data consisted of 3000 successive measurements while

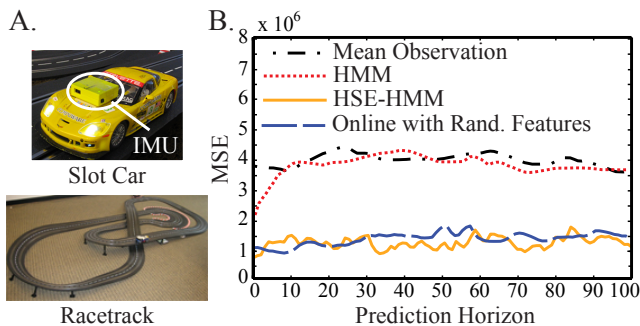


Figure 2: Slot car inertial measurement data. (A) The slot car platform: the car and IMU (top) and the racetrack (bottom). (B) Squared error for prediction with different estimated models. Dash-dot shows the baseline of simply predicting the mean measurement on all frames.

the slot car circled the track controlled by a constant policy. The goal was to learn a model of the noisy IMU data, and, after filtering, to predict future readings.

We trained a 20-dimensional HSE-HMM using the algorithm of Song et al., with tests and histories consisting of 150 consecutive observations. We set the bandwidth parameter of the Gaussian RBF kernels with the “median trick,” and the regularization (ridge) parameter was 10^{-4} . For details see Song et al. (2010).

Next we trained a 20-dimensional TPSR with random Fourier features to *approximate* the Gaussian RBF kernel. We generated 25000 features for the tests and histories and 400 features for current observations, and then used the online spectral algorithm to learn a model. Finally, to provide some context, we learned a 20-state discrete HMM (with 400 levels of discretization for observations) using the Baum-Welch EM algorithm run until convergence.

For each model, we performed filtering for different extents $t_1 = 100, 101, \dots, 250$, then predicted an image which was a further $t_2 = 1, 2, \dots, 100$ steps in the future. The squared error of this prediction in the IMU’s measurement space was recorded, and averaged over all the different filtering extents t_1 to obtain means which are plotted in Figure 2(B).

The results demonstrate that the online spectral learning algorithm with a large number of random Fourier features does an excellent job matching the performance of the HSE-HMM, and suggest that the online spectral learning algorithm is a viable alternative to HSE-HMMs when the amount of training data grows large.

Modeling Video

Next we look at the problem of mapping from video: we collected a sequence of 11,000 160×120 grayscale frames at 24 fps in an indoor environment (a camera circling a conference room, occasionally switching directions; each full circuit took about 400 frames). This data was collected by hand, so the camera’s trajectory is quite noisy. The high frame rate and complexity of the video mean that learning an

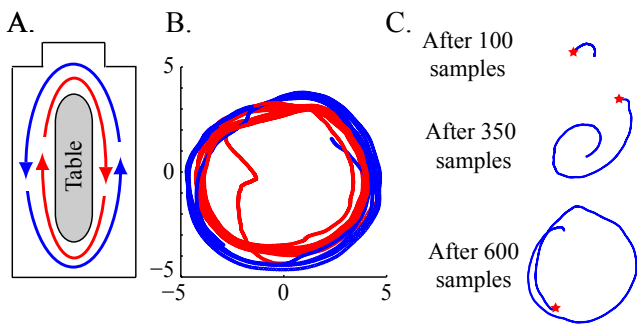


Figure 3: Modeling video. (A.) Schematic of the camera’s environment. (B.) The second and third dimension of the learned belief space (the first dimension contains normalization information). Points are colored red when the camera is traveling clockwise and blue when traveling counterclockwise. The learned state space separates into two manifolds, one for each direction, connected at points where the camera changes direction. (The manifolds appear on top of one another, but are separated in the fourth latent dimension.) (C.) Loop closing: estimated historical camera positions after 100, 350, and 600 steps. Red star indicates current camera position. The camera loops around the table, and the learned map “snaps” to the correct topology when the camera passes its initial position.

accurate model requires a very large dataset. Unfortunately, a dataset of this magnitude makes learning an HSE-HMM difficult or impossible: e.g., the similar but less complex example of Song et al. used only 1500 frames.

Instead, we used random Fourier features and an online TPSR to approximate a HSE-HMM with Gaussian RBF kernels. We used tests and histories based on 400 sequential frames from the video, generated 100,000 random features, and learned a 50-dimensional TPSR. To duplicate this setup, the batch TPSR algorithm would have to find the SVD of a $100,000 \times 100,000$ matrix; by contrast, we can efficiently update our parameters by incorporating 100,000-element feature vectors one at a time and maintaining 50×50 and $50 \times 100,000$ matrices.

Figure 3 shows our results. The final learned model does a surprisingly good job at capturing the major features of this environment, including both the continuous location of the camera and the discrete direction of motion (either clockwise or counterclockwise). Furthermore, the fact that a general-purpose online algorithm learns these manifolds is a powerful result: we are essentially performing simultaneous localization and mapping in a difficult loop closing scenario, *without* any prior knowledge (even, say, that the environment is three-dimensional, or whether the sensor is a camera, a laser rangefinder, or something else).

Conclusion

We presented spectral learning algorithms for TPSR models of partially-observable nonlinear dynamical systems. In particular, we showed how to *update* the parameters of a TPSR

given new batches of data, and built on these updates to develop an efficient *online* spectral learning algorithm. We also showed how to use random projections in conjunction with TPSRs to efficiently approximate HSE-HMMs. The combination of random projections and online updates allows us to take advantage of powerful Hilbert space embeddings while handling training data sets that are orders of magnitude larger than previous methods, and therefore, to learn models that are too complex for previous methods.

Acknowledgements

Byron Boots and Geoffrey J. Gordon were supported by ONR MURI grant number N00014-09-1-1052. Byron Boots was supported by the NSF under grant number EEE-0540865. We would additionally like to thank Sajid Siddiqi and David Wingate for helpful discussions.

References

Boots, B., and Gordon, G. 2010. Predictive state temporal difference learning. In Lafferty, J.; Williams, C. K. I.; Shawe-Taylor, J.; Zemel, R.; and Culotta, A., eds., *Advances in Neural Information Processing Systems 23*. 271–279.

Boots, B.; Siddiqi, S. M.; and Gordon, G. J. 2010. Closing the learning-planning loop with predictive state representations. In *Proceedings of Robotics: Science and Systems VI*.

Bowling, M.; McCracken, P.; James, M.; Neufeld, J.; and Wilkinson, D. 2006. Learning predictive state representations using non-blind policies. In *Proc. ICML*.

Brand, M. 2006. Fast low-rank modifications of the thin singular value decomposition. *Linear Algebra and its Applications* 415(1):20–30.

Hsu, D.; Kakade, S.; and Zhang, T. 2009. A spectral algorithm for learning hidden Markov models. In *COLT*.

Jaeger, H. 2000. Observable operator models for discrete stochastic time series. *Neural Computation* 12:1371–1398.

Katayama, T. 2005. *Subspace Methods for System Identification*. Springer-Verlag.

Littman, M.; Sutton, R.; and Singh, S. 2002. Predictive representations of state. In *Advances in Neural Information Processing Systems (NIPS)*.

Pearl, J. 2000. *Causality: models, reasoning, and inference*. Cambridge University Press.

Rabiner, L. R. 1989. A tutorial on hidden Markov models and selected applications in speech recognition. *Proc. IEEE* 77(2):257–285.

Rahimi, A., and Recht, B. 2007. Random features for large-scale kernel machines. In *Neural Information Processing Systems*.

Rosencrantz, M.; Gordon, G. J.; and Thrun, S. 2004. Learning low dimensional predictive representations. In *Proc. ICML*.

Siddiqi, S.; Boots, B.; and Gordon, G. J. 2010. Reduced-rank hidden Markov models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS-2010)*.

Singh, S.; James, M.; and Rudary, M. 2004. Predictive state representations: A new theory for modeling dynamical systems. In *Proc. UAI*.

Soatto, S., and Chiuso, A. 2001. Dynamic data factorization. Technical Report UCLA-CSD 010001, UCLA.

Sondik, E. J. 1971. The optimal control of partially observable Markov processes. PhD. Thesis, Stanford University.

Song, L.; Boots, B.; Siddiqi, S. M.; Gordon, G. J.; and Smola, A. J. 2010. Hilbert space embeddings of hidden Markov models. In *Proc. 27th Intl. Conf. on Machine Learning (ICML)*.

Van Overschee, P., and De Moor, B. 1996. *Subspace Identification for Linear Systems: Theory, Implementation, Applications*. Kluwer.