

# Learning Multi-Step Predictive State Representations

**Lucas Langer**  
McGill University  
Canada

**Borja Balle**  
Lancaster University  
United Kingdom

**Doina Precup**  
McGill University  
Canada

## Abstract

Recent years have seen the development of efficient and provably correct spectral algorithms for learning models of partially observable environments arising in many applications. But despite the high hopes raised by this new class of algorithms, their practical impact is still below expectations. One reason for this is the difficulty in adapting spectral methods to exploit structural constraints about different target environments which can be known beforehand. A natural structure intrinsic to many dynamical systems is a multi-resolution behaviour where interesting phenomena occur at different time scales during the evolution of the system. In this paper we introduce the multi-step predictive state representation (M-PSR) and an associated learning algorithm that finds and leverages frequent patterns of observations at multiple scales in dynamical systems with discrete observations. We perform experiments on robot exploration tasks in a wide variety of environments and conclude that the use of M-PSRs improves over the classical PSR for varying amounts of data, environment sizes, and number of observations symbols.

## 1 Introduction

Learning models of partially observable dynamical systems is a problem that arises in many practical applications, including robotics, medical monitoring, market analysis, etc. Predictive state representations (PSRs) [Littman *et al.*, 2001; Singh *et al.*, 2004; Rosencrantz *et al.*, 2004], provide a theoretically justified approach in which at any time-step the current state is modeled as a set of predictions about the future evolution of the system, conditioned on past outcomes. The most popular of these models are linear PSRs, in which predictions for any future trajectories can be obtained as a linear combination of a finite set of “core” predictions. A variety of spectral algorithms for learning linear PSRs have been proposed in recent works, e.g. [Boots *et al.*, 2011; Hamilton *et al.*, 2013].

Spectral algorithms are appealing because of their strong theoretical properties, such as statistical consistency and

learning efficiency [Hsu *et al.*, 2009; Bailly *et al.*, 2010]. Unfortunately, their practical uptake is still limited, in part, due to the fact that these general-purpose algorithms are not designed to leverage the structural regularities frequently found in applications. This is a challenging question that needs to be tackled in order to facilitate more efficient spectral learning algorithms for specific applications. Recent work along these lines has focused on sparse states spaces [Hamilton *et al.*, 2013] and special structures encountered in epigenomics [Zhang *et al.*, 2015]. In this paper, we focus on a specific type of structure which arises frequently in dynamical systems: behaviour which takes place at different time scales. This type of structure is often leveraged to provide efficient algorithms in signal processing. Similar gains can be obtained by using multiple time scales in planning and reinforcement learning, e.g. [Sutton *et al.*, 1999].

Our approach is based on a new class of PSRs which we call the multi-step PSR (M-PSR). Like the classical linear PSR, M-PSRs represent the dynamics of a partially observable system with discrete observations. However, unlike PSRs, M-PSRs are able to capture structure occurring at multiple time scales by representing transitions between states over multi-step observations. Our main result is a spectral learning algorithm for M-PSRs combined with a data-driven strategy for multi-step observation selection. Through an extensive empirical evaluation, we show that in environments where characteristic multi-step observations occur frequently, M-PSRs improve the quality of learning with respect to classical PSRs. This improvement is uniform over a range of environment sizes, number of observation symbols, and amounts of training data.

## 2 The Multi-Step PSR

A linear *predictive state representation* (PSR) for an autonomous dynamical system with discrete observations is a tuple  $\mathcal{A} = \langle \Sigma, \alpha_e, \alpha_\infty, \{\mathbf{A}_\sigma\}_{\sigma \in \Sigma} \rangle$  where:  $\Sigma$  is a finite set of possible observations,  $\alpha_e, \alpha_\infty \in \mathbb{R}^n$  are vectors of initial and final weights, and  $\mathbf{A}_\sigma \in \mathbb{R}^{n \times n}$  are the transition operators associated with each possible observation. The dimension  $n$  is the number of states of  $\mathcal{A}$ . Formally, a PSR is a *weighted finite automata* (WFA) [Droste and Vogler, 2009] computing a function given by the probability distribution of sequences of observations in a partially observable dynamical system with a finite number of states. The function  $f_{\mathcal{A}} : \Sigma^* \rightarrow \mathbb{R}$  com-

puted by  $\mathcal{A}$  on input  $x = x_1 \cdots x_t$  is given by:

$$f_{\mathcal{A}}(x_1 \cdots x_t) = \alpha_{\epsilon}^{\top} \mathbf{A}_{x_1} \cdots \mathbf{A}_{x_t} \alpha_{\infty} = \alpha_{\epsilon}^{\top} \mathbf{A}_x \alpha_{\infty}.$$

The value of  $f_{\mathcal{A}}(x)$  is interpreted as the probability that the system produces the sequence of observations  $x = x_1 \cdots x_t$  starting from the initial state specified by  $\alpha_{\epsilon}$ . Depending on the semantics of the model, this can be a probability that the system generates  $x$  and *stops*, or the probability that the system generates  $x$  and *continues*. Given a partial history  $u$  of the evolution of the system the state of a PSR can be updated from  $\alpha_{\epsilon}$  to  $\alpha_u = \alpha_{\epsilon} \mathbf{A}_u$ . This update allows for conditional queries about future observations. For example, the probability of observing a string  $v$  given that we have already observed  $u$  is  $f_{\mathcal{A},u}(v) = \alpha_u \mathbf{A}_v \alpha_{\infty} / \nu_{\mathcal{A}}(u)$ , where  $\nu_{\mathcal{A}}(u)$  is a normalizing constant.

To define a multi-step PSR, we augment a PSR with two extra objects: a set of *multi-step observations*  $\Sigma' \subset \Sigma^+$  containing non-empty strings formed by basic observations, and a *coding function*  $\kappa : \Sigma^* \rightarrow \Sigma'^*$  that takes a string of basic observations and produces an equivalent string composed of multi-step observations. The choice of  $\Sigma'$  and  $\kappa$  can be customized for each application, and will typically capture frequent patterns of observations arising in a particular environment. For the sake of simplicity and to avoid degenerate situations, we assume these objects satisfy the following requirements:

1. The set  $\Sigma'$  must contain all symbols in  $\Sigma$ ; i.e.  $\Sigma \subseteq \Sigma'$
2. The function  $\kappa$  satisfies  $\partial(\kappa(x)) = x$  for all  $x \in \Sigma^*$ , where  $\partial : \Sigma'^* \rightarrow \Sigma^*$  is the *decoding morphism* between free monoids given by  $\partial(z) = z \in \Sigma^*$  for all  $z \in \Sigma'$ . Note this implies that  $\kappa(\epsilon) = \epsilon$ ,  $\kappa(\sigma) = \sigma$  for all  $\sigma \in \Sigma$ , and  $\kappa$  is injective.

Using this notation, we define a *multi-step PSR* (M-PSR) as a tuple  $\mathcal{A}' = \langle \Sigma, \Sigma', \kappa, \alpha_{\epsilon}, \alpha_{\infty}, \{\mathbf{A}_{\sigma}\}_{\sigma \in \Sigma'} \rangle$  containing a PSR with observations in  $\Sigma'$ , together with the basic observations  $\Sigma$  and the corresponding coding function  $\kappa$ . In addition to the standard PSR function  $f_{\mathcal{A}'} : \Sigma'^* \rightarrow \mathbb{R}$ , an M-PSR also defines the function  $f'_{\mathcal{A}'} : \Sigma^* \rightarrow \mathbb{R}$  given by  $f'_{\mathcal{A}'}(x) = f_{\mathcal{A}'}(\kappa(x))$ . In many cases we will abuse notation and write  $f_{\mathcal{A}'}$  for  $f'_{\mathcal{A}'}$  when there is no risk of confusion.

## 2.1 Examples of M-PSRs

We now describe several examples of M-PSR, putting special emphasis on models that will be used in our experiments.

### Base M-PSR

A PSR with a single observation  $\Sigma = \{\sigma\}$  can be used to measure the time – i.e. number of discrete time-steps – until a certain event happens [Bacon *et al.*, 2015]. In this case, a natural approach to build an M-PSR for timing models is to build a set of multi-step observations containing sequences of  $\sigma$ 's whose lengths are powers of a fixed base. That is, given an integer  $b \geq 2$ , the set of multi-step observations is  $\Sigma' = \{\sigma, \sigma^b, \sigma^{b^2}, \dots, \sigma^{b^K}\}$  for some positive  $K$ . A natural choice of coding map in this case is the one that represents any length  $t \geq 0$  as a number in base  $b$ , with the difference that the largest power  $b$  that is allowed is  $b^K$ . This corresponds to writing (in a unique way)  $t = t_0 b^0 + t_1 b^1 + t_2 b^2 +$

$\dots + t_K b^K$ , where  $0 \leq t_k \leq b - 1$  for  $0 \leq k \leq K$ , and  $t_K \geq 0$ . With this decomposition, the coding map is given by  $\kappa(\sigma^t) = (\sigma^{b^K})^{t_K} (\sigma^{b^{K-1}})^{t_{K-1}} \dots (\sigma^b)^{t_1} (\sigma)^{t_0}$ . Note that we choose to write powers of longer multi-step observations first, followed by powers of shorter multi-step observations. For further reference, we will call this model the *Base M-PSR*.

Base M-PSRs can also be extended to the case with multiple basic observations,  $|\Sigma| > 1$ . For example, if there are two observations  $\Sigma = \{\sigma_1, \sigma_2\}$ , we can take  $\Sigma' = \{\sigma_1, \sigma_2, \sigma_1^b, \sigma_2^b, \sigma_1^{b^2}, \sigma_2^{b^2}, \dots, \sigma_1^{b^K}, \sigma_2^{b^K}\}$ . The encoding map  $\kappa$  first splits the string into sequences of consecutive repeated symbols and then uses the same encoding as before. For example:  $\kappa(\sigma_1^5 \sigma_2^3) = (\sigma_1^2)^2 (\sigma_1) (\sigma_2^2) (\sigma_2)$  when using  $b = 2$  and  $K = 1$ .

### Tree M-PSR

A more flexible generalization of the Base M-PSR for the multiple observation case is what we call the *Tree M-PSR*. In a Tree M-PSR, we set  $\Sigma' = \{x \in \Sigma^* : |x| \leq L\}$ , where  $L$  is a parameter. The decoding map  $\kappa$  first splits a string  $x$  as  $x = u_1 u_2 \cdots u_m u_f$ , where  $|u_i| = L$  for  $1 \leq i \leq m$  and  $|u_f| = |x| - (m \cdot L) < L$ , and then sets  $\kappa(x) = (u_1)(u_2) \cdots (u_m)(u_f)$ . Note that this encoding promotes the use of multi-step symbols corresponding to longer strings more often than those corresponding to shorter strings. Note that we have  $|\Sigma'| = O(|\Sigma|^L)$ , so in practice  $L$  must remain small if we want the M-PSR to be representable using a small number of parameters.

### Data-Driven M-PSR

The constructions above have some parameters that can be tuned depending on the particular application, but in general they are not directly dependent on the sequences of observations more frequently produced by the target environment. Intuition suggests that the performance of M-PSRs will depend heavily on how well  $\Sigma'$  reflects frequent observation patterns; this is corroborated by the experiments in Section 4. Thus, we develop an algorithm for adaptively choosing  $\Sigma'$  from data. For this to work, we need to provide a generic coding function  $\kappa$  that can be applied to any M-PSR. This encoding needs to deliver good predictive performance and computationally cheap queries. We call the output of this learning process a *Data-Driven M-PSR*. The details of this approach are described in the next section.

## 3 Learning Algorithm for M-PSR

This section describes a learning algorithm for M-PSR extending standard spectral techniques [Boots *et al.*, 2011] and combining them with an algorithm for building an extended set of symbols  $\Sigma'$  containing frequent patterns of observations.

### 3.1 Spectral Learning Algorithm

We start by extending the spectral learning algorithm to M-PSR under the assumption that  $\kappa$  and  $\Sigma'$  are known. In this case, the learning procedure only needs to recover the operators  $\mathbf{A}_{\sigma}$  for all  $\sigma \in \Sigma'$ , and the initial and final weights  $\alpha_{\epsilon}, \alpha_{\infty}$ . We first recall some basic notation about Hankel

matrices [Carlyle and Paz, 1971; Fliess, 1974], and then proceed to describe the learning algorithm. To simplify the description of the learning algorithm we assume that the function  $f : \Sigma^* \rightarrow \mathbb{R}$  associated with the target M-PSR can be evaluated for every string. In practice these are unknown, but can be effectively estimated from data because they correspond to probabilities of observations.

Given  $f : \Sigma^* \rightarrow \mathbb{R}$ , we will use its *Hankel matrix* representation  $\mathbf{H}_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ , which is an infinite matrix with rows and columns indexed by strings in  $\Sigma^*$  and whose entries satisfy  $\mathbf{H}_f(u, v) = f(uv)$ . To efficiently work with this matrix we only consider finite sub-blocks indexed by a finite set of prefixes  $\mathcal{P} \subset \Sigma^*$  and suffixes  $\mathcal{S} \subset \Sigma^*$ . Both  $\mathcal{P}$  and  $\mathcal{S}$  are input parameters given to the algorithm; see [Balle *et al.*, 2012] for a discussion of general strategies to choose these in practice. The pair  $\mathcal{B} = (\mathcal{P}, \mathcal{S})$  is called a basis, and it determines a sub-block  $\mathbf{H}_{\mathcal{B}} \in \mathbb{R}^{\mathcal{P} \times \mathcal{S}}$  of  $\mathbf{H}_f$  with entries given by  $\mathbf{H}_{\mathcal{B}}(u, v) = \mathbf{H}_f(u, v)$  for all  $u \in \mathcal{P}$  and  $v \in \mathcal{S}$ . For a fixed basis, we also consider the vectors  $\mathbf{h}_{\mathcal{S}} \in \mathbb{R}^{\mathcal{S}}$  with entries given by  $\mathbf{h}_{\mathcal{S}}(v) = \mathbf{H}_f(\epsilon, v)$  for every  $v \in \mathcal{S}$ , and  $\mathbf{h}_{\mathcal{P}} \in \mathbb{R}^{\mathcal{P}}$  with  $\mathbf{h}_{\mathcal{P}}(u) = \mathbf{H}_f(u, \epsilon)$ .

Note that the definitions above only depend on  $\Sigma$ . In order to recover operators  $\mathbf{A}_{\sigma}$  for all  $\sigma \in \Sigma'$  we need to consider multi-step shifts of the finite Hankel matrix  $\mathbf{H}_{\mathcal{B}}$ . In particular, given  $\sigma \in \Sigma'$  we define the sub-block  $\mathbf{H}_{\sigma} \in \mathbb{R}^{\mathcal{P} \times \mathcal{S}}$  whose entries are given by  $\mathbf{H}_{\sigma}(u, v) = f(u\sigma v)$ . This can be interpreted either using the lift  $f(\kappa(u)\sigma\kappa(v))$  or the decoding  $f(u\partial(\sigma)v)$ , but the actual value  $\mathbf{H}_{\sigma}$  will be the same.

Using the notation above we can give a simple description of the learning algorithm. Suppose the basis  $\mathcal{B}$  and the desired number of states  $n$  are given. The algorithm starts by collecting a set of sampled trajectories and uses them to estimate the matrices  $\mathbf{H}_{\mathcal{B}}, \mathbf{H}_{\sigma} \in \mathbb{R}^{\mathcal{P} \times \mathcal{S}}$  and vectors  $\mathbf{h}_{\mathcal{P}} \in \mathbb{R}^{\mathcal{P}}, \mathbf{h}_{\mathcal{S}} \in \mathbb{R}^{\mathcal{S}}$ . Then, it takes the truncated SVD  $\mathbf{U}_n \mathbf{D}_n \mathbf{V}_n^{\top}$  of  $\mathbf{H}_{\mathcal{B}}$ , where  $\mathbf{D}_n \in \mathbb{R}^{n \times n}$  contains the first  $n$  singular values of  $\mathbf{H}_{\mathcal{B}}$ , and  $\mathbf{U}_n \in \mathbb{R}^{\mathcal{P} \times n}$  and  $\mathbf{V}_n \in \mathbb{R}^{\mathcal{S} \times n}$  contain the first left and right singular vectors respectively. Finally, it computes the transition operators of the M-PSR as  $\mathbf{A}_{\sigma} = \mathbf{D}_n^{-1} \mathbf{U}_n^{\top} \mathbf{H}_{\sigma} \mathbf{V}_n$  for all  $\sigma \in \Sigma'$ , and the initial and final weights as  $\alpha_{\epsilon}^{\top} = \mathbf{h}_{\mathcal{S}}^{\top} \mathbf{V}_n$  and  $\alpha_{\infty} = \mathbf{D}_n^{-1} \mathbf{U}_n^{\top} \mathbf{h}_{\mathcal{P}}$ . This algorithm yields an M-PSR with  $n$  states. It was proved in [Boots *et al.*, 2011] that this algorithm is statistically consistent for PSRs (under a mild condition on  $\mathcal{B}$ ). A straightforward argument shows that these same guarantees extend to M-PSR.

### 3.2 A Generic Coding Function

Given  $\Sigma$  and  $\Sigma'$ , a generic coding function  $\kappa : \Sigma^* \rightarrow \Sigma'^*$  can be obtained by minimizing the coding length  $|\kappa(x)|$  of every string  $x \in \Sigma^*$ . More formally, consider the coding

$$\kappa(x) = \operatorname{argmin}_{z \in \Sigma', x=yz, |y| < |x|} |\kappa(y)(z)|$$

Note that for the single observation case, this is equivalent to the optimal coin change problem, which is a textbook example of dynamic programming. This has the advantage of minimizing the number of operators  $\mathbf{A}_{\sigma}$  that will need to be multiplied to compute the value of the M-PSR on a string  $x$ . At the same time, operators expressing long transition sequences can capture the contributions of all intermediate states even if

the chosen model size is too small to represent every state traversed by a single-observation model. We will see in Section 4 that this is one of the reasons why M-PSRs show better performance than PSR for smaller models.

Algorithm 1 gives pseudocode for computing this function. It inductively computes the optimal string encoding for the prefix  $x_1 \cdots x_i$  for all  $1 \leq i \leq |x|$ . This can be obtained by minimizing over all  $\sigma \in \Sigma'$  which terminate at the index  $i$  of  $x$ . We use the following notation:

*bestEnc*: A map from indices  $i$  of the query string  $x$  to the optimal encoding of  $x[:i]$ .

*minEnc*: A map from indices  $i$  of the query string  $x$  to  $|\text{bestEncoding}[i]|$

*opsEnd*: A map from indices  $i$  of  $x$  to the set of strings in  $\Sigma'$ :  $\{s \in \Sigma', x[i - |s| : i] = s\}$

---

#### Algorithm 1 Encoding Algorithm

---

**INPUT:**  $x$

**OUTPUT:**  $\kappa(x)$

```

1: procedure DPENCODE
2:    $\text{bestEnc}[0] = x[0]$ 
3:    $\text{minEnc}[0] = 0$ 
4:   for  $i$  in  $[1, |x|]$  do
5:      $\text{opsEnd}[i] \leftarrow \{s \in \Sigma' : x[i - |s| : i] = s\}$ 
6:   end for
7:   for  $i$  in  $[1, |x|]$  do
8:      $\text{bestOp} \leftarrow \text{null}$ 
9:      $m \leftarrow 0$ 
10:    for  $s \in \text{opsEnd}[i]$  do
11:       $t \leftarrow \text{minEnc}[i - |s|] + 1$ 
12:      if  $\text{bestOp} = \text{null}$  or  $t < m$  then
13:         $m \leftarrow t$ 
14:         $\text{bestOp} \leftarrow s$ 
15:      end if
16:    end for
17:     $\text{minEnc}[i + 1] \leftarrow m$ 
18:     $\text{bestEnc}[i + 1] \leftarrow \text{bestEnc}[i - |\text{bestOp}|] + \text{bestOp}$ 
19:  end for
20:  return  $\text{bestEnc}[|x|]$ 
21: end procedure

```

---

### 3.3 State Update

When using classical PSR in online environments one typically updates the state vector every time a new observation is available. This eliminates the need for repeatedly transforming the initial state over the whole history of observations  $h$  when making predictions. In the case of M-PSRs, the dynamic programming algorithm for minimizing the length of string encodings provides a naturally convenient way to perform efficient state updates. To do so, we cache past state vectors along with their encoding length. When a new observation  $\sigma$  is read, we determine the encoding for the extended observation string  $h\sigma$  with the same recurrence relation as in the previous section. Because this minimization is over  $\{s \in \Sigma' : \exists p \in \Sigma^* ps = h\sigma\}$ , one needs to cache at most  $\max_{s \in \Sigma'} |s|$  state vectors and encoding lengths.

### 3.4 Greedy Selection of Multi-Step Observations

---

**Algorithm 2** Base Selection Algorithm

---

**INPUT:**  $Train, Sub_M$

**OUTPUT:**  $\Sigma'$

```

1: procedure BASE SELECTION
2:    $\Sigma' \leftarrow \Sigma$ 
3:    $bestEnc \leftarrow null$ 
4:   for all  $obs \in Train$  do
5:      $bestEnc[obs] \leftarrow |obs|$ 
6:   end for
7:    $i \leftarrow 0$ 
8:   while  $i < numOps$  do
9:      $bestOp \leftarrow null$ 
10:     $m \leftarrow 0$ 
11:    for all  $s \in Sub_M$  do
12:       $c \leftarrow 0$ 
13:      for all  $obs \in Train$  do
14:         $c \leftarrow c + DPEnc(obs, \Sigma' \cup s) -$ 
           $bestEnc[obs]$ 
15:      end for
16:      if  $c < m$  then
17:         $bestOp \leftarrow s$ 
18:         $m \leftarrow c$ 
19:      end if
20:    end for
21:     $\Sigma' \leftarrow \Sigma' \cup bestOp$ 
22:    for all  $obs \in Train$  do
23:       $bestEnc[obs] \leftarrow DPEnc(obs, \Sigma')$ 
24:    end for
25:     $i \leftarrow i + 1$ 
26:  end while
27:  return  $\Sigma'$ 
28: end procedure

```

---

Selecting multi-step transition sequences to build  $\Sigma'$  can be achieved with an adaptive greedy algorithm, depicted in Algorithm 2. A  $\Sigma'$  that correctly captures the frequent observations produced by a target environment should promote short encodings when coupled with the coding function  $\kappa$  described above. In practice, we want  $\Sigma'$  to contain substrings appearing in the training data which are long, frequent, and diverse. From an intuitive standpoint, one can view structure in observation sequences as relating to the level of entropy in the distribution over multi-step observations produced by the system, and interpret this selection method as a data compression scheme.

In general terms, the algorithm works as follows. In a pre-processing step, we find all possible substrings in  $\Sigma^*$  that appear in the training dataset. For computational reasons this can be constrained only to the  $M$  most frequent substrings, where  $M$  is a parameter chosen by the user. Frequency of occurrence is measured by number of training trajectories that contain a given substring. The construction of  $\Sigma'$  is then initialised by  $\Sigma$  and proceeds iteratively by adding a new multi-step symbol at each phase. A phase starts by evaluating all substrings in terms of how much reduction in the number of transition operators used to encode the whole training set

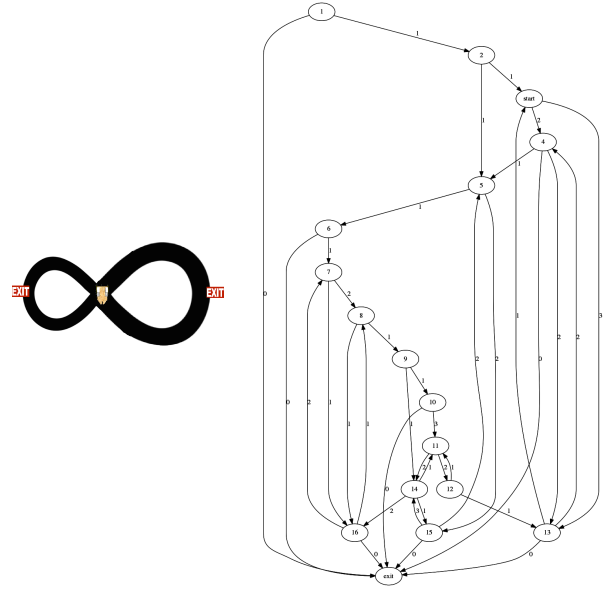


Figure 1: The Double Loop environment (left) and the Pacman Labyrinth (right)

would be achieved if the symbol was added to  $\Sigma'$ . The algorithm then adds to  $\Sigma'$  the multi-step symbol corresponding to the best substring, i.e. the one that would reduce the most the whole coding cost (with respect to  $\kappa$ ). More formally, at the  $i$ th iteration the algorithm finds:

$$\operatorname{argmin}_{u \in sub_M} \sum_{x \in train} |\kappa_{\Sigma'_i \cup \{u\}}(x)|,$$

where  $train$  is the training set,  $sub_M$  is the set of substrings under consideration of length at most  $M$ ,  $\Sigma'_i$  is the set of multi-step observations at the beginning of phase  $i$  and we use  $\kappa_{\Sigma'}$  to denote the encoding function with respect to a given set of multi-step observations for clarity. The algorithm terminates after  $\Sigma'$  reaches a predetermined size.

## 4 Experiments

In this section, we assess the performance of PSRs and M-PSRs in various environments, with different model sizes, and learned from both large and small datasets. For all the plots, the x-axis is model size of the PSR/M-PSRs and the y-axis is an error measurement of the learned PSR/M-PSRs.

In all the experiments, an agent is positioned in a partially observable environment and navigates stochastically based on state to state transition probabilities. An observation symbol is produced on every transition. When the agent exits, we record the concatenation of the symbols produced, which represents the observation sequence. We perform experiments in two environments: a Double Loop maze and a Pacman-style environment, depicted in Figure 1.

For the Base M-PSR (i.e., the timing models), we construct the empirical Hankel matrix by taking  $\mathcal{P}, \mathcal{S} = \{\sigma^i, \forall i \leq N\}$ , where  $N$  is an application-dependent parameter. For Double Loop environments, we set  $N = 150$ , while for the Pacman domain,  $N = 600$ . For these choices of  $N$ , we verify that as the amount of data increases, the learned PSR with the true

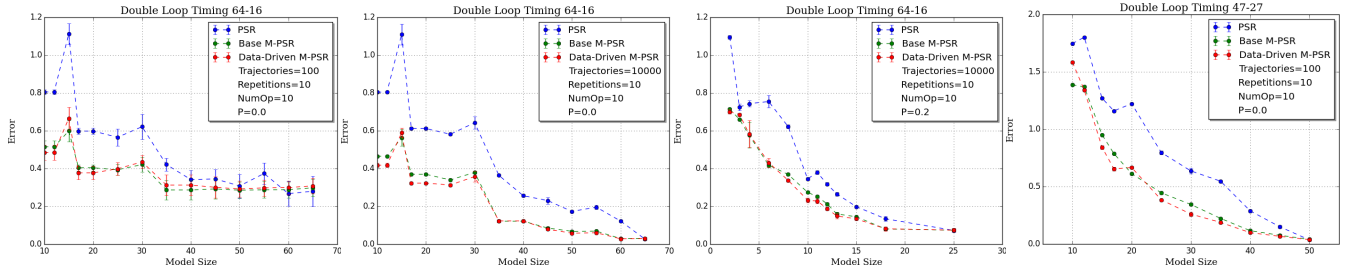


Figure 2: Results for the Double Loop of size 64-16: (1) 100 samples,  $p = 0$  (2) 10000 samples,  $p = 0$  (3) 10000 samples,  $p = 0.2$ ; and the Double Loop of size 47-27: (4) 100 samples,  $p = 0$ .

model size converges to the true model. For the Base M-PSR, we set  $b = 2$ ,  $K = 8$ , so that the longest string in  $\Sigma'$  is  $\sigma^{256}$ .

For the tasks with multiple observations, a slightly more complex approach is required to choose  $\mathcal{P}$  and  $\mathcal{S}$ . For the prefixes  $\mathcal{P}$ , we select the  $k$  most frequent prefixes from our observation set. For the suffixes  $\mathcal{S}$ , we take all suffixes that occur in  $\mathcal{P}$ . We also require prefix completeness. That is, if  $p'$  is a prefix of  $p \in \mathcal{P}$ , then  $p' \in \mathcal{P}$ . This heuristic for constructing empirical Hankel matrices was also used in [Balle *et al.*, 2012]. For the Base M-PSR, we take  $K = 8$  and  $B = 2$  for both symbols  $\Sigma = \{g, b\}$ , where  $g$  stands for green and  $b$  for blue. For the Tree M-PSR, we set  $L = 7$  for a total of 128 operators, a far larger limit than for the other M-PSRs.

To measure the performance of our learning algorithms we compute the difference between the values computed by  $f$ , the true probability distribution over observations, and  $\hat{f}$ , the function associated with the learned M-PSR/PSR. In our environments,  $f$  can be computed exactly, as we have access to the underlying HMMs. Since there are infinitely many strings, we only compute the difference over a fixed set  $T \subset \Sigma^*$  which depends on the problem. Thus, reported errors are given by

$$\|f - \hat{f}\|_T = \sqrt{\sum_{x \in T} (f(x) - \hat{f}(x))^2}.$$

For the Base M-PSR (timing) case, we take  $T = \{\sigma^i : i \leq C\}$  with  $C = 600$ . Importantly,  $C$  has to be large enough to ensure the error measure provides a measure of accuracy over the most frequent strings, say  $\sum_{i=0}^C f(\sigma^i) > 0.99$ . For the multiple observation case, we take all possible strings that can be produced from the prefixes and suffixes in our dataset:  $T = \{ps : p \in \mathcal{P}, s \in \mathcal{S}\}$ . In all our experiments we use the same  $\mathcal{P}$  and  $\mathcal{S}$  to learn PSR and the different types of M-PSR.

#### 4.1 Double Loop Timing

We start by considering the Double Loop environment, and the task of learning timing models. The length of a loop is the number of states in the loop. A trajectory begins with the agent at the intersection of the two loops. Here, the agent has a 50% probability of entering either loop. At intermediate states, the agent moves to the next state in the loop with probability  $1 - p$  or remains in its current state with probability  $p$ . Exit states are located halfway into each loop. The agent leaves the environment at exit states with 50% probability.

Figure 2 (1-3) provides results for the case in which one loop consists of 64 states and the other of 16 states. Panel

(1) presents the results when learning from 100 observation sequences, while the results in panel (2) are obtained using 10000 observations. In both cases, the M-PSRs outperform the simple PSR for small model sizes. Panel (3) shows the effect of varying the self-transition probability to  $p = 0.2$ , in order to simulate noise in the environment. We find that the environment with noisy loops is more compressible, that is, one can achieve better performance for low model sizes, but the performance becomes worse as the model size attains the environment's true size. M-PSRs still significantly outperform the standard PSR for reduced model sizes. Panel (2) in Figure 4 illustrates the effect of  $K$  on the performance of the Base M-PSRs. We note that larger values of  $K$  have better performance, up to  $K = 7$ .

In Figure 4 (2), we illustrate how varying the number of multi-step transition operators affects performance for the 64-16 double loop. As expected, a higher number of operators improves performance, but the effect tapers off at about 20 operators. Here, the most important operators are:  $\{\sigma, \sigma^8, \sigma^{24}, \sigma^{32}, \sigma^{72}\}$  which are closely tied to the environment's periodic structure.

In Figure 2 (4), we plot the results of a 47-27 labyrinth. Here, because of the lengths of the loops, observations will not be as compactly expressed from the Base M-PSR. Again, M-PSRs outperform the standard PSR for reduced model sizes.

#### 4.2 Pacman Labyrinth Timing

We now turn our attention to the environment on the right of Figure 1. The weights in the transitions represent the “length” of the corridors they represent, and at each state the agent has equal probability of following any corridor starting in that state. To test the effect of the geometry of the environment we add a parameter  $sF$  (stretch factor) which is used to multiply the length of the corridors (by adding extra intermediate states). In the left panels of Figure 3 we vary the number of observations used for learning. M-PSRs outperform the traditional PSR regardless of the amount of data. In the right two panels, we vary the stretch factor parameter, while keeping the size of the dataset fixed. We find that a higher value of  $sF$  provides increased improvement of the M-PSR relative to the standard PSR.

#### 4.3 Multiple Observations: Coloured Loops

To test the case of multiple observations, we construct a Double Loop environment where the first loop has length  $l_1 = 27$

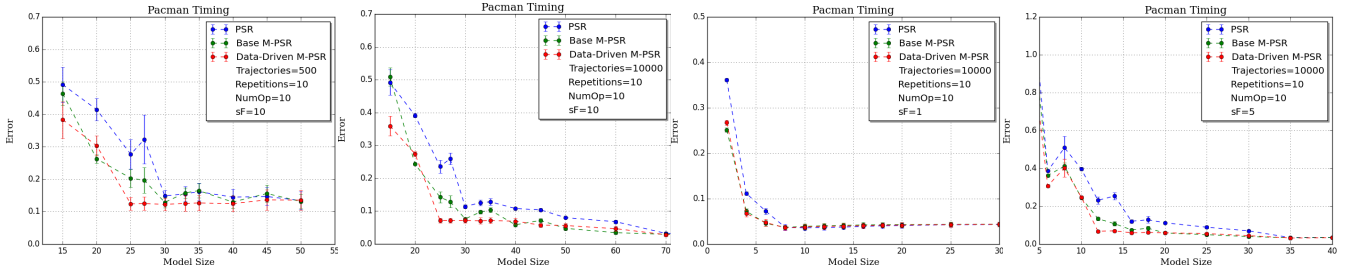


Figure 3: Results for the Pacman environment with (in order from left to right) low amount of data, high amount of data and stretch factor equal to 1 and 5 (right panels)

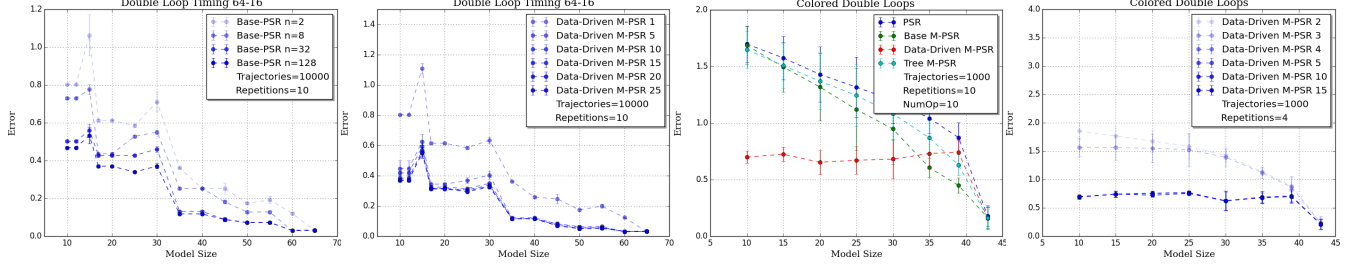


Figure 4: (1) Effect of exponent in Base M-PSR (Double Loop 64-16) (2) Effect of the number of operators (Double Loop 64-16) (3) Effect of model size (Colored Double Loop 27-17) (4) Effect of the number of operators (Colored Double Loop 27-17)

and observations are green. The second loop is blue, with length  $l_2 = 17$ . We fix the length of each trajectory at  $3(l_1 + l_2)$ . We build empirical estimates for the Hankel matrix as follows:

$$f(x) = \frac{|\text{train} \cap x\Sigma^*|}{|\text{train} \cap \Sigma^{\geq |x|}|}.$$

This means that the PSRs will compute the probability of  $x$  occurring as a prefix.

As for the timing case, we find that M-PSRs perform far better, especially the Data-Driven M-PS. Results are displayed in panel (3) of Figure 4. In the panel (4) we vary the number of multi-step transition operators learned. The important operators learned are  $\{g, b, g^{27}, b^{17}\}$ , which is again very encouraging, as it reflects the structure of this particular environment.

## 5 Discussion

We presented an approach to leveraging multiple temporal scales in the behaviour of dynamical systems, in order to learn predictive state representations. The proposed model, M-PSR, can be learned by adapting existing spectral approaches. To the best of our knowledge, the only other work that attempts to include temporal abstraction in PSRs is due to [Wolfe and Singh, 2006], who use temporally extended actions, or options [Sutton *et al.*, 1999], and build PSRs on top of these. However, this model has a very different flavour from our approach, as we bundle together observations into multiple steps. In particular, our approach is applicable even when there are no actions, such as in the case of HMMs, whereas this previous work requires structure in the action

space. In all the experiments we conducted, M-PSRs offer significantly better predictive performance for reduced model sizes than PSRs. In addition, Data-Driven PSRs offer improvements over generic M-PSRs by learning the transition operators specific to the environment, which is very important when prior information about appropriate time scales for modelling is not known. Our evaluation focused on illustrating the advantage of the proposed model especially when the amount of data available is small, and in noisy environments. These experiments suggest that M-PSR with small number of states provide much better approximations to a given environment than PSR with the same number of states. We think this is a very promising property, and we plan to study it in detail in future work. We used simulated domains in our experiments in order to compare with a ground truth. However, we anticipate that the proposed models would be useful in real tasks, for example in financial market prediction. In this case, one could discretize market levels into bins, and use M-PSRs to learn predictive models over multiple time steps, in order to capture both long-term trends and fast fluctuations. We also plan to experiment with this approach in the domain of processing physiological signal recordings, which similarly exhibit short and long-term trends. On the theoretical side, it would be interesting to analyze the existence of an “optimal” set of symbols for the data-driven M-PSR, and to develop further algorithms for finding it.

## References

[Bacon *et al.*, 2015] P.-L. Bacon, B. Balle, and D. Precup. Learning and planning with timing information in Markov Decision Processes. In *UAI*, 2015.

- [Bailly *et al.*, 2010] R. Bailly, A. Habrard, and F. Denis. A spectral approach for probabilistic grammatical inference on trees. In *ALT*, 2010.
- [Balle *et al.*, 2012] B. Balle, A. Quattoni, and X. Carreras. Local loss optimization in operator models: A new insight into spectral learning. *International Conference on Machine Learning (ICML)*, 2012.
- [Boots *et al.*, 2011] B. Boots, S. Siddiqi, and G. Gordon. Closing the learning planning loop with predictive state representations. *International Journal of Robotic Research*, 2011.
- [Carlyle and Paz, 1971] J. W. Carlyle and A. Paz. Realizations by stochastic finite automata. *Journal of Computer Systems Science*, 1971.
- [Droste and Vogler, 2009] W. Droste, M. Kuich and H. Vogler. *Handbook of weighted automata*. Springer, 2009.
- [Fliess, 1974] M. Fliess. Matrices de Hankel. *Journal de Mathématiques Pures et Appliquées*, 1974.
- [Hamilton *et al.*, 2013] William L. Hamilton, Mahdi M. Fard, and Joelle Pineau. Modelling sparse dynamical systems with compressed predictive state representations. In *ICML*, 2013.
- [Hsu *et al.*, 2009] D. Hsu, S.M. Kakade, and T. Zhang. A spectral algorithm for learning Hidden Markov Models. In *COLT*, 2009.
- [Littman *et al.*, 2001] M. L. Littman, R. S. Sutton, and S. Singh. Predictive representations of state. *Neural Information Processing Systems Conference (NIPS)*, 2001.
- [Rosencrantz *et al.*, 2004] M. Rosencrantz, G. Gordon, and S. Thrun. Learning low dimensional predictive representations. In *ICML*, 2004.
- [Singh *et al.*, 2004] S. Singh, M. R. James, and M. R. Rudary. Predictive state representations: A new theory for modeling dynamical systems. In *UAI*, 2004.
- [Sutton *et al.*, 1999] R. S Sutton, D. Precup, and S. Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 1999.
- [Wolfe and Singh, 2006] B. Wolfe and S. Singh. Predictive state representations with options. *International Conference on Machine Learning (ICML)*, 2006.
- [Zhang *et al.*, 2015] Chicheng Zhang, Jimin Song, Kamalika Chaudhuri, and Kevin Chen. Spectral learning of large structured hmms for comparative epigenomics. *Neural Information Processing Systems (NIPS)*, 2015.