

## An effective reinforcement learning with automatic construction of basis functions and sequential approximation

Nobuhito Nanjoh, Takeshi Mori

Graduate School of Information Science

Nara Institute of Science and Technology

8916-5, Takayama, Ikoma, Nara, 630-0192, Japan

Shin Ishii

Graduate School of Informatics

Kyoto University

Gokasho, Uji, Kyoto, 611-0011, Japan

### Abstract

In reinforcement learning, a large state space makes the value function estimation infeasible, and hence it is often approximated as a linear combination of basis functions whose linear coefficients constitute the parameter. However, the basis function construction needs some prior knowledge and is often troublesome. To overcome this difficulty, Keller et al. proposed an automatic basis function construction technique [2], but it may cause serious computational cost. We propose a novel approach to this context, where the computational cost can be drastically reduced from that of the existing method, due to our sequential approximation.

### 1 Introduction

Reinforcement learning (RL) is a machine learning method through interactions with environments to maximize the long-term reward accumulation (return). In usual RL schemes, it is important to estimate the value function which predicts the return starting from each state. In many realistic RL problems, however, a large state space makes the value function estimation in its original function space infeasible, and hence, the value function approximation using a parametric linear model comes to be important; the value function is represented as a linear combination of basis functions whose linear coefficients constitute the parameter. While the parameter is obtained by the conventional RL methods based on sample data, the basis functions must be determined by a designer at hand in advance, but the determination is very difficult due to the small amount of available prior knowledge. A poor setting of basis functions may increase the potential error of approximating the value function, and a large approximation error makes the policy's control fail.

One possible idea to deal with this problem is to construct the basis functions automatically based on the available data in an online manner. Keller et al. proposed an automatic basis function construction technique [2], in which appropriate basis functions were generated based on tem-

poral difference (TD) error. However, this approach had some difficulties. First, the basis functions were generated in each time step, and then, all the basis functions, whose number could be huge, were used to approximate the value function by means of the least-squares optimization. This may cause serious computational cost, as the number of bases increases. Second, the value function approximation can be unstable due to the singularity in the least squares solution, especially as the number of highly correlated basis functions increases.

To overcome these problems, in this study, we propose a novel approach to the automatic construction of basis functions, where the number of basis functions and the computational cost can be drastically reduced from those of the existing method. According to our method, the basis functions, are constructed in to approximate the TD error at each time step instead of the value function, and the learned parameters are fixed in the subsequent optimizations. In other words, our method decomposes whole the approximation problem in the Markov decision processes (MDPs) into small approximation problems which are easily to be solved with a small number of basis functions, so that the parameter is optimized within each sub-problem. This method has some theoretical foundations, given by Bertsekas and Castanon [3], which showed that the sequential calculation of analytic approximation of the TD error with appropriate basis functions can geometrically reduce the approximation error along the whole time-series. Computer simulation shows our method could drastically reduce not only the number of basis functions but also the computational time in comparison to the Keller et al.'s method [2].

### 2 MDPs and the value function approximation

We consider finite MDPs, which is composed of the discrete time  $t$ , a set of states  $\mathcal{S}$ , a set of actions  $\mathcal{A}$ , a set of transition probabilities  $\mathcal{P}$ , and a scalar reward  $r \in \mathcal{R}$ . At time  $t$ , the agent selects an action  $a_t \in \mathcal{A}$  according to a

stationary policy  $\pi$  at a state  $s_t \in \mathcal{S}$ , and then it moves to the next state  $s_{t+1} \in \mathcal{S}$  and simultaneously receives the reward  $r_{t+1}$ . The objective is to find a policy that maximizes the value function:

$$V(s) = E \left[ \sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s \right] \quad (1)$$

where  $\gamma \in [0, 1)$  is a discount factor.

One of the approaches of searching for the optimal policy is policy iteration [1]. The policy iteration is constructed of two steps, i.e., the policy evaluation step and the policy improvement step. In the policy evaluation step, the value function  $V$  for the current policy  $\pi$  is calculated or approximated. In the policy improvement step, the policy  $\pi$  is improved by using the learned value function. In this article, we focus on the policy evaluation step. The Bellman equation under  $\pi$  is defined as

$$V(s) = E [r_{t+1} + \gamma V(s_{t+1}) | s_t = s]. \quad (2)$$

For the sake of later convenience, we introduce the matrix notations; the value function vector  $\mathbf{V} \in \mathcal{R}^{|\mathcal{S}|}$  whose  $i$ -th element is the value at the  $i$ -th state in  $\mathcal{S}$ , the state transition matrix  $\mathbf{P} \in \mathcal{R}^{|\mathcal{S}| \times |\mathcal{S}|}$  whose  $i, j$ -th element is the state transition probability under policy  $\pi$  from the  $i$ -th state to the  $j$ -th state in  $\mathcal{S}$ , and the reward vector  $\mathbf{r} \in \mathcal{R}^{|\mathcal{S}|}$  whose  $i$ -th element is the expected reward with respect to the next state conditioned on the  $i$ -th state and policy  $\pi$ . Then, the Bellman equation Eq.(2) is represented as  $\mathbf{V} = \mathbf{r} + \gamma \mathbf{P}\mathbf{V}$ .

In the dynamic programming context, the policy evaluation in average operates successive approximations of the value function  $\mathbf{V}$  using Bellman operator  $T(\cdot)$ :

$$\mathbf{V}^{k+1} = T(\mathbf{V}^k) \equiv \mathbf{r} + \gamma \mathbf{P}\mathbf{V}^k, \quad (3)$$

where  $k$  denotes an iteration step and  $\mathbf{V}^k$  converges to  $\mathbf{V}$  in the limit of  $k \rightarrow \infty$  [1]. In usual RL settings, the vector  $\mathbf{r}$  and the matrix  $\mathbf{P}$  are unknown, so the value function is estimated by using sample trajectory,  $\langle s_0, r_1, s_1, r_2, \dots, s_t, r_{t+1}, \dots \rangle$ .

In realistic RL problems, however, the state space is large, and the value function estimation is very difficult. Then the value function is often approximated by using a parametric linear model, that is, represented as a linear combination of basis function whose linear coefficients constitute the parameter:

$$V(s) \approx \sum_{m=1}^M \phi_m(s)' \theta_m = (s)' \boldsymbol{\theta}, \quad (4)$$

where  $(\cdot)'$  is the transpose, and the basis function vector  $\phi(s)$  and the parameter  $\boldsymbol{\theta}$  are both  $M$ -dimensional vectors. Note that the designer of the system must prepare

the basis functions at hand. The parameter  $\boldsymbol{\theta}$  is tuned in the policy evaluation step by using the sample trajectory. One of the learning methods is least-squares TD learning [4], which is the closed-form estimation method for the linearly-approximated value function.

### 3 Previous method for automatic basis function construction and aggregation theory

In usual, the basis function construction requires designer's trial and error. The poor setting of basis functions may increase the approximation error of the value function, and a large approximation error makes the policy's update fail. One possible idea to deal with this problem is to construct the basis functions automatically based on the available data in an online manner.

Keller et al. [2] proposed an automatic basis function construction technique, in which appropriate basis functions were generated based on the aggregation theory, proposed formerly by Bertsekas and Castanon [3]. According to this theory, automatically produced aggregated states are used to speed up the value iteration in dynamic programming, even if the original operations of the Bellman operator (Eq.(3)) are very slow. In section 5, we show that our proposal method is also based on this theory, actually a more appropriate application of this theory rather than Keller et al.'s. In this theory, the Bellman operator (Eq.(3)) is replaced by aggregation iterations of the form:

$$\mathbf{V}^{k+1} = \mathbf{V}^k + \boldsymbol{\gamma}, \quad (5)$$

where  $\boldsymbol{\gamma} \in \{0, 1\}^{|\mathcal{S}| \times M}$  is such an aggregation matrix that  $\gamma_{ij} = 1$  if state  $i$  is assigned to group  $j$  or 0 otherwise. The grouping of the states is determined by the ranking of the Bellman residual:  $T(\mathbf{V}^k) - \mathbf{V}^k$ .  $\boldsymbol{\gamma}$  is the parameter vector corresponding to  $\boldsymbol{\gamma}$ , which is derived by analytical calculation.  $\boldsymbol{\gamma}$  is the orthogonal projected vector of the Bellman residual onto the subspace spanned by sparse matrix  $\boldsymbol{\gamma}$ . By using the Bellman operator  $T$ ,

$$T(\mathbf{V}^{k+1}) = T(\mathbf{V}^k) + \boldsymbol{\gamma}, \quad (6)$$

is obtained, and by using the definition of orthogonal projection:  $\Pi = (\boldsymbol{\gamma}' \boldsymbol{\gamma})^{-1} \boldsymbol{\gamma}'$ , the following holds [3]:

$$T(\mathbf{V}^{k+1}) - \mathbf{V}^{k+1} = \mathbf{g}_1 + \mathbf{g}_2, \quad (7)$$

where

$$\begin{aligned} \mathbf{g}_1 &= (\mathbf{I} - \Pi) (T(\mathbf{V}^k) - \mathbf{V}^k), \\ \mathbf{g}_2 &= (\mathbf{I} - \Pi) \boldsymbol{\gamma}. \end{aligned}$$

Based on the above ranking scheme of the Bellman residual  $T(\mathbf{V}^k) - \mathbf{V}^k$ , the maximum norm of the first approximation error  $\mathbf{g}_1$  can be minimized under a restricted number

of groups  $M$ . Moreover, it is known experimentally that  $g_2$  often decreases as  $g_1$  decreases [3]. So we may focus only on the minimization of  $g_1$ .

In the Keller et al.'s method, by applying this theory to an RL context, the value function is updated as

$$\begin{aligned} \bar{V}_{k+1} &= \Phi_k \theta_k + \theta_{\text{new}} \\ &= [\Phi_k \quad ] \begin{matrix} \theta_k \\ \theta_{\text{new}} \end{matrix} = \Phi_{k+1} \theta_{k+1}, \end{aligned} \quad (8)$$

where  $\Phi_k \in \{0, 1\}^{|\mathcal{S}| \times kM}$  is the basis function matrix:  $\Phi_k = [ \phi_k(1), \phi_k(2), \dots, \phi_k(|\mathcal{S}|) ]'$ ,  $\phi_k(s) \in \{0, 1\}^{kM}$  is a  $kM$ -dimensional binary vector corresponding to a state  $s$ , and  $M$  is the number of basis functions produced in each step. In this method, the basis functions are constructed iteratively based on the TD error (the Bellman residual for a sampled state), and the parameters are evaluated by least-squares TD learning (LSTD) [4]. By using a sample trajectory, the parameter vector  $\theta_{k+1}$  is optimized so as to minimize the cost function:

$$J_{\text{LSTD}}(\theta_{k+1}) = \frac{1}{2t} \sum_{i=0}^{t-1} (r_{i+1} - V^k(s_{i+1}) - \phi_{k+1}(s_i)' \theta_{k+1})^2, \quad (9)$$

where  $t$  is the current time step. Then, in each iteration step  $k$ ,  $kM$ -dimensional vector  $\theta_k$  is estimated by the least-squares optimization.

Because the number of the parameters increases by  $M$  after each iteration, the computation (in the order of  $\mathcal{O}((kM)^3)$  due to the inverse calculation of the design matrix) can be intractable. Furthermore, the value function approximation can be unstable due to the singularity of the design matrix as the number of highly correlated basis functions increases.

#### 4 Sequential approximation method for automatic basis function construction

To overcome these difficulties, we propose a novel approach for automatic basis function construction, where the computational cost is only of  $\mathcal{O}(M^3)$ , which is independent of  $k$ , and stable learning can be achieved by avoiding the singularity. While in the previous approach shown above, the value function is approximated by all parameters, we consider the approximation of the TD error instead of the value function with restricted parameters  $\theta_{\text{new}}$ , and learned parameters  $\theta_k$  being fixed after each optimization and embedded in  $V^k$ . To this end, we adopt least-squares policy evaluation learning (LSPE) [5] instead of LSTD. In this method, the TD error instead of the value function is approximated by least-squares optimization, so as to have a form:

$$\bar{V}_{k+1} = V_k + \theta_{\text{new}}. \quad (10)$$

Note that only  $M$ -dimensional vector  $\theta_{\text{new}}$  is learned, and the others are fixed. The cost function is given by

$$J(\theta_{\text{new}}) = \frac{1}{2t} \sum_{i=0}^{t-1} (r_{i+1} + V^k(s_{i+1}) - V^k(s_i) - \psi(s_i)' \theta_{\text{new}})^2, \quad (11)$$

where  $\psi(s_i) \in \{0, 1\}^M$  is a newly produced basis function and is the transpose of the row vector in (8), evaluated at the state  $s_i$ . TD error:  $r_{i+1} + V^k(s_{i+1}) - V^k(s_i)$  is the target of  $\psi(s_i)' \theta_{\text{new}}$  in this regression problem. In the minimization of this objective function with respect to  $\theta_{\text{new}}$ , the least-squares solution is given as

$$\theta_{\text{new}} = B^{-1} c, \quad (12)$$

where

$$\begin{aligned} B &= \frac{1}{t} \sum_{i=0}^{t-1} \psi(s_i) \psi(s_i)', \\ c &= \frac{1}{t} \sum_{i=0}^{t-1} \psi(s_i) (r_{i+1} + V^k(s_{i+1}) - V^k(s_i))'. \end{aligned}$$

Then, the TD error can be successively approximated by sequentially constructed basis functions. Note that the newly produced basis functions are only used for each sequential approximation. Then, the computational cost is of  $\mathcal{O}(M^3)$ , which is much smaller than  $\mathcal{O}((kM)^3)$  of the Keller et al.'s method. Moreover, it can be seen that our method is the straightforward application of the Bertsekas and Castanon's aggregation theory shown in Eq.(5)-(7), since the TD error is orthogonally projected onto the space  $\Psi$ . The algorithm is depicted in Fig 1.

Figure 1: *single iteration step*

##### input

$h_t = \langle s_0, r_1, \dots, s_t, r_t \rangle$ : a sample trajectory until  $t$   
 $M$ : the number of basis functions  
 $V_k$ : estimated value function

##### For $i = 1:t-1$

$$e_i \leftarrow r_{i+1} + V_k(s_{i+1}) - V_k(s_i)$$

##### End

$$\Psi \leftarrow \text{GenerateBasis}(\{e\}, M)$$

$$\theta_{\text{new}} \leftarrow \text{LSPE}(h_t, V_k, \Psi)$$

$$V_{k+1} \leftarrow V_k + \Psi \theta_{\text{new}}$$

##### return $V_{k+1}$

The process of the adaptive basis function construction is described by *GenerateBasis*. This function generates basis functions according to the TD errors  $e$ . The parameter  $\theta_{\text{new}}$  is estimated by LSPE (Eq.(12)), and then the value function is updated. Since the previous parameters are fixed, serious instability due to the singularity can be avoided in our method.

## 5 Experimental results

To verify the performance of our method, we evaluate the maximum norm<sup>1</sup> between the true value function and the estimated value functions by several methods in simulation experiments. Since the maximum norm is associated with the worst case policy improvement, this measure is important for the convergence property of RL [1]. We compared our method with the following three methods:

- LSPE with fixed basis functions, which is a usual setting of conventional RL methods,
- LEPE with randomly generated basis functions,
- LSTD with adaptive construction of basis functions (Keller et al.'s method).

By comparing with (a), we can show the effectiveness of sequential construction of basis functions. By comparing with (b), we can show the effectiveness of the adaptive constructions based on TD error. By comparing with (c), we can show the effectiveness of our method compared with the Keller et al.'s method. We assume that the number of states is 100, the discount factor is 0.97, the number of basis functions ( $M$ ) is 2, and the transitional matrix under the fixed policy  $\pi$  is as  $P = Z(0.03 \times P_{\text{rand}} + (1 - 0.03) \times I)$ , where the components of  $P_{\text{rand}}$  are drawn from the uniform distribution in  $(0, 1]^{|S| \times |S|}$ ,  $I$  is a unit matrix, and  $Z$  is a normalization term. Then the true value function is given by  $V_{\text{true}} = (I - P)^{-1} r$ . Samples are drawn from the transition matrix  $P$  step by step starting from an initial state  $s_0$ . We generated new basis functions and the estimated value function every 10 sample generations. Fig 2 shows the performance comparison where the horizontal axis is the number of generated samples. The performances were averaged over 100 independent runs.

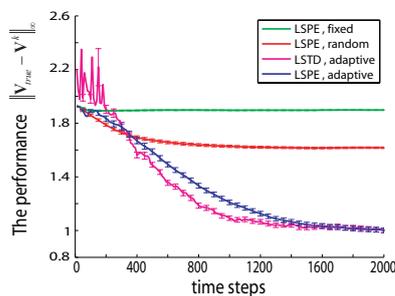


Figure 2: The performance comparisons

With setting (a), the performance saturated immediately due to the restricted expression of fixed basis functions. With setting (b), the performance gradually increased, actually much better than (a) every time step. The Keller et

<sup>1</sup>For an  $N$ -dimensional vector  $\mathbf{x} = [x_1, x_2, \dots, x_N]$ , the maximum norm is defined as  $\|\mathbf{x}\|_{\infty} = \max(x_1, x_2, \dots, x_N)$ .

al.'s method (setting (c)) drastically improved the performance, but it showed substantial instability behavior due to simultaneous optimization of all the parameters. In the following learning phase, the design matrix was always singular, so we had to use the singular value decomposition. Moreover, the computational cost of  $\mathcal{O}((kM)^3)$  was a hazard for examination of longer time-steps behaviors. In our method, the performance was improved similar to (c), but did not show any instability in contrast to (c). Furthermore, since the computational cost is of only  $\mathcal{O}(M^3)$ , the learning can be continued for any period in on-line manner; this is an important feature for scaling-up the application.

## 6 Conclusion and future works

We proposed a novel scheme for automatic construction of basis functions, without huge computational cost. In this method, the value function was estimated by sequential approximation of the TD errors with adaptive basis functions. In our experiments, we showed that the performance of our method was much better than previous methods in spite of small computational cost.

Although our method was formalized by using aggregated discrete states, many realistic problems have continuous state space. We need some modifications of our method to be applicable to continuous-state space problems. One possibility is to extend the value function approximator to that with continuous-state basis functions such as the radial basis functions, but this extension hinder direct application of the aggregation theory, and the performance assurance can be abused. So theoretical foundations are desired in continuous-state settings. Furthermore, the combination of the automatic dimensionality reduction method, such as neighborhood component analysis, explored by Keller et al. [2], could be interesting, but this orientation does not have sufficient theoretical backgrounds yet.

## References

- [1] Bertsekas, D. and Tsitsiklis, J. (1996), *Neuro-Dynamic Programming*, ISBN 1-8865 29-10-8
- [2] Keller, P., Mannor, S., and Precup, D. (2006), Automatic Basis Function Construction for Approximate Dynamic Programming and Reinforcement Learning. *International Conference on Machine Learning*
- [3] Bertsekas, D. and Castanon, D. (1989), Adaptive aggregation methods for infinite horizon dynamic programming. *IEEE Transactions of Automatic Control*, 34,589-598
- [4] Bradtke, S. and Barto, A. (1996), Linear Least-Squares Algorithms for Temporal Difference Learning. *Machine Learning*, 22, 33-57
- [5] Bertsekas, D. and Yu, H. (2006), Convergence results for some temporal difference methods based on least squares. (Technical Report LIDS-2697).