

Unsupervised Basis Function Adaptation for Reinforcement Learning

Edward Barker · Charl Ras

Received: date / Accepted: date

Abstract When using reinforcement learning (RL) algorithms to evaluate a policy it is common, given a large state space, to introduce some form of approximation architecture for the value function (VF). The exact form of this architecture can have a significant effect on the accuracy of the VF estimate, however, and determining a suitable approximation architecture can often be a highly complex task. Consequently there is a large amount of interest in the potential for allowing RL algorithms to adaptively generate (i.e. to learn) approximation architectures.

We investigate a method of adapting approximation architectures which uses feedback regarding the frequency with which an agent has visited certain states to guide which areas of the state space to approximate with greater detail. We introduce an algorithm based upon this idea which adapts a state aggregation approximation architecture on-line.

Assuming S states, we demonstrate theoretically that – provided the following relatively non-restrictive assumptions are satisfied: (a) the number of cells X in the state aggregation architecture is of order $\sqrt{S} \ln S \log_2 S$ or greater, (b) the policy and transition function are close to deterministic, and (c) the prior for the transition function is uniformly distributed – our algorithm can guarantee, assuming we use an appropriate scoring function to measure VF error, error which is arbitrarily close to zero as S becomes large. It is able to do this despite having only $O(X \log_2 S)$ space complexity (and negligible time complexity). We conclude by generating a set of empirical results which support the theoretical results.

Keywords reinforcement learning · unsupervised learning · basis function adaptation · state aggregation

E. Barker and C. Ras
School of Mathematics and Statistics
University of Melbourne, Victoria 3010, Australia
Tel.: +61-3-9035-8877
E-mail: ebarker@student.unimelb.edu.au

1 Introduction

Most commonly used reinforcement learning (RL) algorithms store an estimate of what’s known as the value function (VF). The VF corresponds to a particular policy, and is a mapping from each state-action pair to a real value which reflects the amount of reward the agent will obtain starting from that state-action pair and following the policy in question (Sutton and Barto 1998). In order for an RL algorithm to perform well (i.e. to achieve a high reward), it is important that the VF estimate is as accurate as possible, since it is this estimate which governs how the algorithm will update its policy.

Traditional RL algorithms such as $TD(\lambda)$ or Q -learning can generate exact VF estimates when dealing with small state and action spaces. However, when environments are complex (with large state or action spaces), applying such algorithms directly becomes too computationally demanding. As a result it is common to introduce some form of architecture with which to approximate the VF, for example a parametrised set of functions (Sutton and Barto 1998; Bertsekas and Tsitsiklis 1996). One issue when introducing VF approximation, however, is that the accuracy of the algorithm’s VF estimate is highly dependent upon the exact form of the architecture chosen (it may be, for example, that no element of the chosen set of parametrised functions closely fits the VF).

Accordingly, a number of authors have explored the possibility of allowing the approximation architecture to be *learned* by the agent, rather than pre-set manually by the designer – see Buşoniu et al (2009) for an overview. What we might hope to achieve by employing such an approach is to create an RL algorithm which still has relatively low computational demands, but at the same time has increased flexibility, allowing us to apply the algorithm to a wider set of problems without needing to invest time to suitably adapt it in each case. If we assume that the approximation architecture being adapted is linear (so that the VF is represented as a weighted sum of basis functions) such methods are known as *basis function adaptation*.

A simple and perhaps, as yet, under-explored method of basis function adaptation involves using an estimate of the frequency with which an agent has visited certain states to determine which states to more accurately represent. Such methods are “unsupervised” in the sense that no direct reference to the reward or to any estimate of the VF is made. The concept of using visit frequencies in an unsupervised manner is not completely new (Menache et al 2005; Bernstein and Shimkin 2010) however it remains relatively unexplored compared to methods which seek to measure the error in the VF estimate explicitly and to then use this error as feedback (Munos and Moore 2002; Bertsekas and Yu 2009; Di Castro and Mannor 2010; Mahadevan et al 2013).

As we will demonstrate, however, unsupervised methods have some distinct advantages: (a) estimates of visit frequencies are cheap to calculate and to store, (b) accurate estimates of visit frequencies can be generated with fewer samples than accurate estimates of, for example, temporal differences (which are a common form of feedback used to adapt basis functions), and (c) under suitable conditions, and applying an appropriate scoring function, such methods can in fact generate very accurate VF estimates, guaranteeing in certain cases scores arbitrarily close to zero. Our overarching objective in this article is to more closely examine unsupervised methods and seek to quantify, where possible, these advantages.

It is point (c) which is perhaps the most surprising, and it forms the substance of the article’s main result. For any fixed policy you will have a stationary distribution describing the probability of being in each state. Suppose (i) the policy and transition function for a given problem are close to deterministic, and (ii) the prior for the transition function is uniformly distributed (this will be more precisely described below). We will show that, under these conditions, an agent which follows an arbitrary policy will, on average, tend to spend almost all of its time in a small subset of the state space. Indeed, if the state space is of size S , there is a theoretical upper bound on the average size of this subset: $O(\sqrt{S} \ln S)$. Provided we have enough basis functions to individually represent the states in this subset, unsupervised basis function adaptation methods can ensure that the VF is arbitrarily well estimated over this subset. If the scoring function we apply is weighted by the probability of visiting each state, we can then guarantee an arbitrarily low score. The implication of this is that, under these circumstances, unsupervised methods will perform *at least as well* as other more complex methods, but, compared to these other methods, will do so *more cheaply* (in the sense of sampling required and also computational demands). Whilst conditions (i) and (ii) encompass many important and general problems, we will also explore the potential to generalise condition (ii) to encompass a larger set of possible priors.

We also explore these ideas experimentally. Our experimental results suggest that our techniques provide a powerful advantage in many real world settings. Over a set of realistic parameter settings the techniques (when compared to fixed state aggregation) can reduce VF error by an amount in the range of 40-70%. In some cases the experimental results also suggest that the assumptions required by the theory can be relaxed.

As noted above, unsupervised techniques at present are relatively unexplored. Menache et al (2005) provided a brief evaluation of an unsupervised algorithm (to provide a comparison with two more complex adaptation algorithms) in the setting of policy evaluation.¹ Bernstein and Shimkin (2009) examined an algorithm where a set of kernels are progressively split (“once-and-for-all”) based on the visit frequency for each kernel. Their algorithm includes policy updates which incorporate knowledge of uncertainty in the VF estimate. The algorithm we propose below works in conjunction with a state aggregation approximation architecture, employing a form of “cell-splitting” to give state space regions more or less resolution. This bears similarities to a number of approaches examined in the literature (the main difference in our approach is the rules under which cells are split or joined). Moore and Atkeson (1995) provide an early algorithm based on updating a state aggregation architecture, whilst Whiteson et al (2007) examined a basis function construction method involving cell splitting.

Our analysis is new both in terms of the details of the unsupervised algorithm we outline (it is designed so as to minimise memory requirements, in particular to ensure that space complexity is roughly of the order of the number of basis functions, whilst still permitting the approximation architecture to be continuously adapted on-line) and in terms of the theoretical concepts we derive.

¹ Their paper actually found the unsupervised method performed unfavourably compared to the alternative approaches they proposed. However the environment they used to test the algorithm does not satisfy our stated assumptions.

In the remainder of this section we outline the formal framework we will be using. In Section 2 we set out the details of our algorithm (PASA, short for “Probabilistic Adaptive State Aggregation”) which performs unsupervised basis function adaptation based on state aggregation. In Section 3 we outline our main theoretical results. Finally, in Section 4 we set out some empirical results designed to both support and extend the results in Section 3.

1.1 Formal framework

We assume that we have an agent which interacts with an environment over a sequence of iterations $t \in \mathbb{N}$. For each t the agent will be in a particular state² s_i ($1 \leq i \leq S$) and will take a particular action a_j ($1 \leq j \leq A$). Each action is taken according to a *policy* π whereby the probability the agent takes action a_j in state s_i is denoted as $\pi(a_j|s_i)$. We are considering the problem of policy evaluation, so we will always assume that the agent’s policy π is fixed (i.e. does not change as a function of t). The *reward function* R is a mapping from each state-action pair (s_i, a_j) to a real number, such that if the agent is in state s_i and takes action a_j , then it will receive a *reward* $R(s_i, a_j)$. The reward function is assumed to be deterministic and bounded, i.e. $|R(s_i, a_j)| < \infty$ for all (i, j) .

The *transition function* P defines how the agent’s state evolves over time. If the agent is in state s_i and takes an action a_j in iteration t , then the probability it will transition to the state $s_{i'}$ in iteration $t+1$ is given by $P(s_{i'}|s_i, a_j)$. Both P and R are taken as unknown, however we assume we are given a prior distribution for both.

The *value function* Q^π , which maps each of the $S \times A$ state-action pairs to a real value, is defined as follows:

$$Q^\pi(s_i, a_j) = \mathbb{E}_\pi \left(\sum_{t=1}^{\infty} \gamma^{t-1} R(s^{(t)}, a^{(t)}) \middle| s^{(1)} = s_i, a^{(1)} = a_j \right) \quad (1)$$

where the expectation is taken over the distributions of P and π (i.e. for a particular instance of P , not over its prior distribution) and where $\gamma \in [0, 1)$ is known as a *discount factor*. We have used superscript brackets to indicate dependency on the iteration t . Initially the VF is unknown, our objective, given some π , is to learn the VF for a particular instance of P .

Assuming we are able to store an explicit real value for each state-action pair, traditional RL algorithms provide a means of estimating Q^π . These estimates will converge to the correct value as t becomes large (Sutton and Barto 1998). In cases where S or A are large, though, it may be impossible to store $S \times A$ real values. Hence, when dealing with such cases, it is common to employ VF approximation. Once we approximate the VF, however, even if the underlying RL algorithm still converges to generate an estimate of the VF (which is not always guaranteed), we can no longer rely on the estimate being arbitrarily close to the true value Q^π .

One form of VF approximation, *parametrised value function approximation*, involves generating an approximation of the VF using a parametrised set of functions. The goal of the RL algorithm then becomes to find a value for the parameters

² The state space is assumed to be discrete. In order for the results in Section 3 to hold in the case of a continuous state space, additional assumptions would need to be introduced.

so that the VF estimate is as near to the true value as possible. The approximate VF is denoted as \hat{Q}_θ , and, assuming we are approximating over the state space only and not the action space, this value is parametrised by a matrix θ of dimension $X \times A$ (where $X \ll S$). Such an approximation architecture is *linear* if \hat{Q}_θ can be expressed in the form $\hat{Q}_\theta(s_i, a_j) = \varphi(s_i, a_j)^T \theta_j$, where θ_j is the j th column of θ and $\varphi(s_i, a_j)$ is a fixed vector of dimension X for each pair (s_i, a_j) . The XA distinct vectors of dimension S given by $(\varphi(s_1, a_j)_k, \varphi(s_2, a_j)_k, \dots, \varphi(s_S, a_j)_k)$ are called *basis functions*. It is common to assume that $\varphi(s_i, a_j) = \varphi(s_i)$ for all j , in which case we have only X distinct basis functions, and $\hat{Q}_\theta(s_i, a_j) = \varphi(s_i)^T \theta_j$.

A *state aggregation* approximation architecture – see, for example, Singh et al (1995) and Whiteson et al (2007) – is a simple linear approximation architecture which we can define as being a mapping F from each state s_i to a *cell* x_k ($1 \leq k \leq X$). Defining the architecture as a “mapping” implies that every state corresponds to exactly one cell. We will denote as \mathcal{X}_k the set of states in the k th cell. Given a state aggregation approximation architecture, the underlying RL algorithm cannot distinguish states in the same cell, and hence $\hat{Q}_\theta(s_i, a_j)$ will be the same for all states which are in the same cell (θ in this context can be interpreted as a set of weights, one given to each cell-action pair).

1.2 Scoring the value function estimate

If we want to design an algorithm to adapt an approximation architecture, we need a means to assess how well it is doing this. A *scoring function* is used to assess the accuracy of a VF estimate (and can also therefore be used to evaluate an algorithm designed to generate a VF estimate). Many basis function adaptation algorithms use a scoring function as a form of feedback to help guide how the basis functions should be updated. In such cases it is important that the score is something which can be measured computationally.

One commonly used score is the squared error in the VF estimate for each state-action, weighted by the probability of each state-action occurring (Menache et al 2005; Bertsekas and Yu 2009; Di Castro and Mannor 2010). We will refer to this as the *mean squared error* (MSE):

$$\text{MSE} := \sum_{i=1}^S \psi_i \sum_{j=1}^A \pi(a_j | s_i) \left(Q^\pi(s_i, a_j) - \hat{Q}_\theta(s_i, a_j) \right)^2 \quad (2)$$

This is where ψ is a vector of the probability of each state given the stationary distribution associated with π (given some fixed policy π , the transition matrix obtained from π and P has a corresponding stationary distribution, provided the transition matrix is irreducible and aperiodic). Note that the true VF Q^π appears in (2). This value, however, is unknown. Therefore, another commonly used scoring function (which, unlike MSE, can be estimated empirically) uses T^π , the *Bellman operator*, to obtain an approximation of the MSE. This scoring function we denote as L (this is a weighted sum of what is sometimes known as the *Bellman error* at each state-action):

$$L := \sum_{i=1}^S \psi_i \sum_{j=1}^A \pi(a_j | s_i) \left(T^\pi \hat{Q}_\theta(s_i, a_j) - \hat{Q}_\theta(s_i, a_j) \right)^2 \quad (3)$$

where:

$$T^\pi \hat{Q}_\theta(s_i, a_j) := R(s_i, a_j) + \gamma \sum_{i'=1}^S \sum_{j'=1}^A P(s_{i'}|s_i, a_j) \pi(a_{j'}|s_{i'}) \hat{Q}_\theta(s_{i'}, a_{j'}) \quad (4)$$

Our results in Section 3 will be stated in relation to MSE and L . It will be crucial to all of our results that the scoring function is weighted by ψ . Two important comments should be made in relation to this.

The first is that, whilst applying such a weighting appears natural, a scoring function does not necessarily need to be weighted by ψ (or an approximation of ψ). There may be circumstances under which a more appropriate measure of the accuracy of a VF estimate would, for example, weight every state equally (the most appropriate measure to use in each situation would depend on a number of complex factors). We acknowledge this as a limitation of the analysis in Section 3.

The second is that, in investigating unsupervised basis function adaptation methods, we are implicitly making a comparison with methods of basis function adaptation which use explicit scores as a source of feedback (these could be called “supervised” methods). If an algorithm uses a scoring function as feedback, then (irrespective of which scoring function is most appropriate) it is best evaluated in terms of how well it minimises *that particular scoring function*. The scoring function L is an attractive choice to provide feedback for supervised methods since, by weighting the error by the probability of visiting a state, it is possible to generate an estimate of the score without knowing the probability of visiting each state. In fact, any feedback based on a scoring function which is *not* weighted by ψ would implicitly require some way of normalising the score for each state. This in turn would require an estimate for ψ , which implies that $O(S)$ distinct values need to be recorded. Hence, if unsupervised methods can perform comparatively well in terms of minimising probability weighted scoring functions, this is of great importance when comparing such methods to supervised alternatives.

In the definitions above both MSE and L are weighted by the probability of each action occurring. We could redefine MSE and L to weight each action equally – when sampling for L this would not be an issue computationally since π is known to the algorithm, i.e. we can simply divide each sample by $\pi(a_j|s_i)$ (acknowledging that samples from rarely chosen actions would contribute more significantly to the variance of the estimate). Our results in Section 3 will extend to L under such an alternative definition.

2 The PASA algorithm

2.1 Underlying concepts

As we noted in the introduction, circumstances exist under which an agent (adopting a random, or in some cases an arbitrary, policy) will, on average, tend to spend almost all of its time in a relatively small subset of the state space. We examine these circumstances more closely in Section 3. The underlying idea of PASA is to make the VF representation as detailed as possible over this relatively small subset (whilst allowing the representation to be coarser over the remainder of the state space).

PASA is designed to function in conjunction with a state aggregation approximation architecture, and accordingly it updates a mapping F over time. Provided this mapping converges to some fixed mapping F^* (we prove below that it will), then if an RL algorithm such as SARSA (which we have used for the experiments in Section 4) is used to update \hat{Q}_θ for the state aggregation architecture associated with F^* , the estimate \hat{Q}_θ will also converge, due to the fact that SARSA, as well as many other RL algorithms, will, with a linear approximation architecture, converge for fixed policies (Bertsekas and Tsitsiklis 1996). We can then assess \hat{Q}_θ using our scoring function.

Suppose, momentarily, we have sufficient computational space to be able to generate an estimate $\hat{\psi}_i$ of each ψ_i (which we could do, for example, by having S weights and using a stochastic approximation algorithm). We could then design our algorithm roughly as follows: (i) start with an initial coarse set of $B < X$ cells, (ii) calculate $u_k = \sum_{i:i \in \mathcal{X}_k} \hat{\psi}_i$ for $1 \leq k \leq B$, (iii) split the cell with the largest value of u_k , (iv) recalculate u_k for the new $B + 1$ cells, and (v) continue in the same fashion until we have X cells. We could rerun this splitting process every iteration, or, if we prefer, at discrete intervals. Provided ψ_i is large for only a small set of states (i.e. the variance of the elements of ψ is high) the resulting set of cells (basis functions) will tend to have a more detailed representation of the VF in areas of the state space with high stationary probability.

This is the essence of how PASA works, except that, to avoid storing $O(S)$ weights (which we’ve implicitly assumed is impossible), we instead measure the probability of visiting the B “base cells” and of visiting $X - B$ additional cells (which may change over time) which are progressively split from these base cells. We can then estimate the probability of visiting each of the X resulting cells by subtracting estimates from one another. The consequence, as will become apparent, is that PASA converges to the same point as the algorithm described in the paragraph above, whilst requiring only $O(X \log_2 S)$ space complexity. The trade off is that it may take a longer amount of time to converge, although in practical terms this difference would appear to be marginal.

Before setting out how PASA works in detail, it is worthwhile highlighting one aspect of our algorithm. Note that many methods examined in the literature involve what has been termed *basis function construction* (Buşoniu et al 2009), where a set of basis functions are determined by an initial process “once-and-for-all”, which occurs prior to the agent beginning to function “as normal” – examples include Munos and Moore (2002) and Whiteson et al (2007). Such methods work, for example, by progressively adding basis functions until some criteria is satisfied. The PASA algorithm falls into an alternative class of methods which have been termed *basis function optimisation* (Buşoniu et al 2009). The assumption with such methods is that there is a fixed number of basis functions which are progressively updated throughout the whole period the agent functions. This approach has the advantage of being more flexible – the basis functions can adapt to policy changes, or indeed to changes in the environment.

2.2 Details of the algorithm

PASA will store a vector ρ of integers of dimension $X - B$, where $B < X$. Suppose we start with a fixed partition of the state space into B cells, indexed from 1

to B , each of which is approximately the same size. Using ρ we can now define a new partition by splitting (as evenly as possible) the ρ_1 th cell in the original partition. We leave one half of the ρ_1 th cell with the index ρ_1 and give the other half the index $B + 1$ (all other indices stay the same). Taking this new partition (consisting of $B + 1$ cells) we can create a further partition by splitting the ρ_2 th cell. Continuing in this fashion we will end up with a partition containing X cells (which gives us the mapping F). We need some additional mechanisms to allow us to update ρ . Denote by $\mathcal{X}_{i,j}$ the set of states in the i th cell of the j th partition (so $0 \leq j \leq X - B$ and $1 \leq i \leq B + j$). The algorithm will store a vector \bar{u} of real values of dimension X (initialised as a vector of zeroes). This will record the approximate frequency with which certain cells have been visited by the agent. We define a new vector \bar{x} of dimension X accordingly:

$$\bar{x}_i^{(t)} = \begin{cases} I_{\{s^{(t)} \in \mathcal{X}_{i,0}\}} & \text{if } 1 \leq i \leq B \\ I_{\{s^{(t)} \in \mathcal{X}_{i,i-B}\}} & \text{if } B < i \leq X \end{cases} \quad (5)$$

where I is the *indicator function* for a logical statement such that $I_A = 1$ if A is true. We can interpret \bar{x} as follows: \bar{x}_i will be equal to 1 if and only if the current state s falls into any cell which has been split from the original cell $\mathcal{X}_{0,i}$ or $\mathcal{X}_{i,i-B}$ (depending on i). The resulting mapping from each state to a vector \bar{x} we denote as \bar{F} (there is a simple mapping from each vector \bar{x} to one of the X final cells: simply take the mapped-to cell as the highest index k such that $\bar{x}_k = 1$). We then update \bar{u} in each iteration as follows (i.e. using a simple stochastic approximation algorithm):

$$\bar{u}_i^{(t+1)} = \bar{u}_i^{(t)} + \eta \left(\bar{x}_i^{(t)} - \bar{u}_i^{(t)} \right) \quad (6)$$

This is where $\eta \in (0, 1]$ is a constant step-size parameter. To update ρ , at certain intervals $\nu \in \mathbb{N}$ the PASA algorithm performs a sequence of $X - B$ operations. A temporary copy of \bar{u} is made, which we call u . We also store an X dimensional boolean vector Σ and set each entry to zero at the start of the sequence. This keeps track of whether a particular cell has only one state, as we don't want the algorithm to try to split singleton cells. At each stage $k \geq 1$ of the sequence we update ρ as follows (for ρ , if multiple indices satisfy the arg max function, we take the lowest index):

$$\rho_k = \begin{cases} j & \text{if } (1 - \Sigma_{\rho_k})u_{\rho_k} < \max\{u_i : i \leq B + k - 1, \Sigma_i = 0\} - \vartheta \\ \rho_k & \text{otherwise} \end{cases} \quad (7)$$

where

$$j = \arg \max_i \{u_i : i \leq B + k - 1, \Sigma_i = 0\} \quad (8)$$

and where $\vartheta > 0$ is a constant designed to ensure that a (typically small) threshold must be exceeded before ρ is adjusted. We also, after ρ , update:

$$\begin{aligned} u_{\rho_k} &\leftarrow u_{\rho_k} - u_{B+k} \\ \Sigma_i &= I_{\{|\mathcal{X}_{i,k}| \leq 1\}} \text{ for } 1 \leq i \leq B + k - 1 \end{aligned} \quad (9)$$

The idea behind each step in the sequence is that the non-singleton cell $\mathcal{X}_{i,k}$ with the highest value u_i (an estimate of visit frequency which is recalculated at each step) will be split. Details of these steps, as well as the overall PASA

process, are outlined in Algorithm 1. Note that the algorithm calls a procedure to SPLIT cells. This procedure simply updates \bar{F} and Σ given the latest value of ρ . It also calls a CONVERT procedure, which converts the mapping \bar{F} to a mapping F . Both of these procedures are computationally very straightforward. A diagram illustrating the main steps is at Figure 1.

Algorithm 1 The PASA algorithm. Called at each iteration t . Assumes \bar{u} , \bar{F} , F and ρ are stored. Return is void.

```

1: function PASA( $t, s, \eta, \vartheta, \nu$ )
2:    $\bar{x} \leftarrow \bar{F}(s)$ 
3:    $\bar{u} \leftarrow \bar{u} + \eta(\bar{x} - \bar{u})$ 
4:   if  $t \bmod \nu = 0$  then
5:      $u \leftarrow \bar{u}$ 
6:     for  $k \in \{1, \dots, X\}$  do
7:        $\Sigma_k \leftarrow 0$ 
8:     end for
9:     for  $k \in \{1, \dots, X - B\}$  do
10:       $u_{\max} \leftarrow \max\{u_i : i \leq B + k - 1, \Sigma_i = 0\}$ 
11:       $i_{\max} \leftarrow \min\{i : i \leq B + k - 1, u_i = u_{\max}, \Sigma_i = 0\}$ 
12:      if  $(1 - \Sigma_{\rho_k})u_{\rho_k} < u_{i_{\max}} - \vartheta$  then
13:         $\rho_k \leftarrow i_{\max}$ 
14:      end if
15:       $u_{\rho_k} \leftarrow u_{\rho_k} - u_{B+k}$ 
16:       $(\bar{F}, \Sigma) \leftarrow \text{SPLIT}(k, \rho, \bar{F})$ 
17:    end for
18:     $F \leftarrow \text{CONVERT}(\bar{F})$ 
19:  end if
20: end function

```

2.3 Some basic properties of the algorithm

PASA requires only a modest increase in computational resources compared to fixed state aggregation. In relation to time complexity, \bar{u} can be updated in parallel with the RL algorithm's update of θ (and the update of \bar{u} would not be expected to have any greater time complexity than the update to θ if using a standard RL algorithm such as SARSA or Q -learning). The vector ρ can be updated at intervals ν (and this update can also be run in parallel). In practice ν can be large because this allows time for \bar{u} to converge. The mapping from state to cell has a very low order of time complexity: $O(\log_2 S)$ for an RL algorithm using PASA compared to a minimum of $O(\log_2 X)$ for X equally sized cells. Hence, PASA involves effectively no increase in time complexity.

PASA does involve additional space complexity with respect to storing the vector \bar{u} : we must store X real values. If we also store F and \bar{F} (as well as u temporarily) the overall space complexity becomes $O(X \log_2 S)$, although F must also be stored for fixed state aggregation. The RL component has space complexity $O(XA)$ (reflecting the $X \times A$ cell-action pairs), so that the introduction of PASA as a pre-processing algorithm will not materially impact the overall space complexity. (Note also that the space complexity of PASA is independent of A .)

Regarding sampling efficiency, two points can be made. The first is that, since visit frequencies do not depend on each individual action, reward or subsequent trajectory, they can be estimated quickly, much more quickly than, for example, temporal differences. The second point arises when we compare PASA to methods based on explicitly estimating the VF error (supervised methods). Once PASA has converged (we argue below that it will), then the estimate \hat{Q}_θ only needs to converge *once*. In contrast, if the estimate \hat{Q}_θ is used to update the basis functions, then a new value of \hat{Q}_θ must, in principle, be generated each time there is an update to the basis functions. This could have serious consequences for the time required for the process to converge, as \hat{Q}_θ may take a very long time to generate an accurate estimate for a particular set of basis functions. This is particularly so if γ is near 1 or when S or A are large.³

³ Moreover, in the special case where a scoring function weights all actions equally, if A is large or π heavily favours certain actions then an accurate estimate of the Bellman error (obtained, for example, using temporal differences) will require an even larger number of samples, since rarely taken actions will have a high amount of variance.

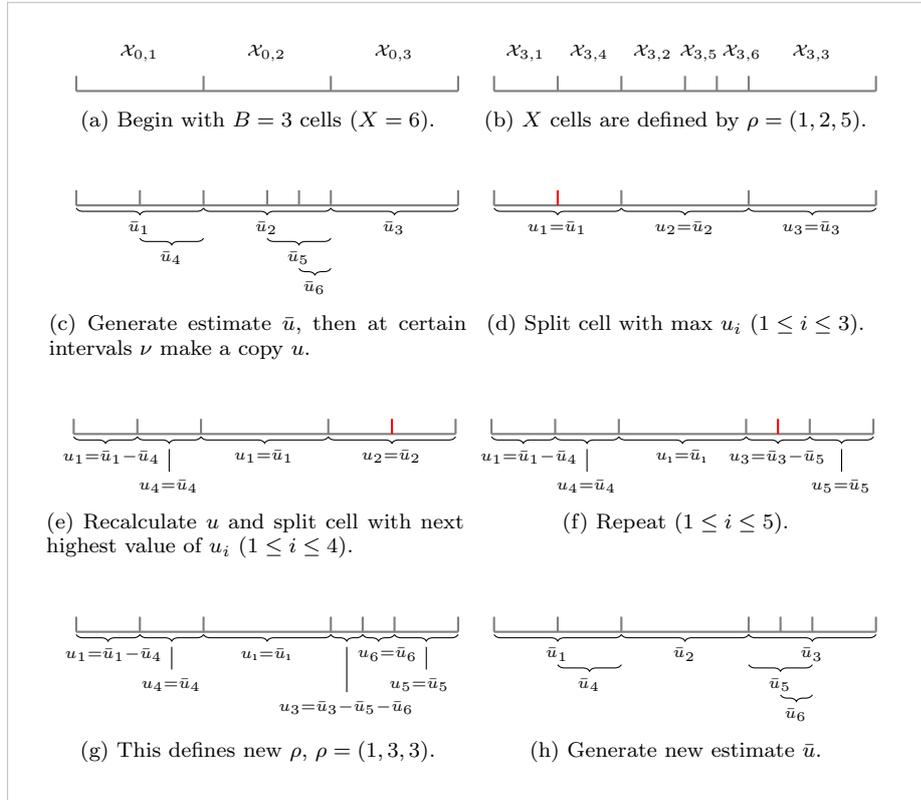


Fig. 1: Simplified summary of PASA algorithm, with $X = 6$ and $B = 3$ (S is arbitrary). The horizontal line represents the set of all states, whilst the short vertical lines represent the breaks between cells.

In relation to the convergence properties of PASA, if we are considering a situation where π is updated, then provided π continues to change PASA will not converge and instead will continuously update F . However, for fixed π , PASA will converge in a particular sense which we now describe. Our outline of PASA assumed a single fixed step-size parameter η . For our proof below it will be easier to suppose that we have a distinct fixed step-size parameter η_k for each element \bar{u}_k of \bar{u} , each of which we can set to a different value (fixed as a function of t). For the remainder of this section η should be understood as referring to this vector of step-size parameters. We will say that some function of t , $x = x^{(t)}$, becomes ε -fixed over τ after T provided T is such that, for all $T' > T$, the value x will remain the same for all t' satisfying $T' \leq t' \leq T' + \tau$ with probability at least $1 - \varepsilon$ (where the probability is taken over the prior distributions for P and R).

Proposition 1 *For every $\tau \in \mathbb{Z}$ and $\varepsilon > 0$ there exists η and T such that the mapping $F^{(t)}$ generated by PASA will become ε -fixed over τ after T .*

Proof Suppose that $\{h_i\}_{i=1}^{X-B}$ is a sequence of real numbers such that $h_i > 0$ for all i and that $\{H_i\}_{i=1}^{X-B}$ is a sequence of closed intervals on the real line such that $\max H_i - \min H_i = h_i$ for all i . Provided we have suitably specified H_i , we can denote:

$$I_i := I_{\{\sum_{j=1}^{B+i-1} \bar{u}_j \in H_i\}} \quad (10)$$

For any $\tau_1, h_1 > 0$ and $\varepsilon_1 > 0$ we can set the elements $\{\eta_k\}_{k=1}^B$ small enough, find H_1 and find some T_1 large enough so that $\sum_{j=1}^B \bar{u}_j$ remains within the interval H_1 for τ_1 iterations with probability at least $1 - \varepsilon_1$ provided T_1 iterations have elapsed. In this way I_1 becomes ε_1 -fixed over τ_1 after T_1 .

This follows by relying on results regarding stochastic approximation algorithms with fixed step-sizes. By replacing each $\bar{u}_k^{(t)}$ with a time scaled equivalent, so that $t \rightarrow \lfloor t/\eta_k \rfloor$, each \bar{u}_k will, as $\eta_k \rightarrow 0$, have an associated deterministic ordinary differential equation (ODE) with a unique solution. Furthermore, the difference between each \bar{u}_k and this ODE converges (again as $\eta_k \rightarrow 0$) weakly in distribution to an Ornstein-Uhlenbeck process, which has a normal stationary distribution with a scaling factor of $\sqrt{\eta_k}$ (Bucklew et al 1993). Hence, due to the fact that each \bar{u}_k is bounded by the unit interval, we can select a scalar η' so that our requirement is satisfied provided $\eta_k \leq \eta'$ for all $1 \leq k \leq B$.

Suppose we have chosen $h_1 < \vartheta/2$. Then ρ_1 will also become ε_1 -fixed over τ_1 after T_1 . This is for the following reason. Suppose each u_i for $1 \leq i \leq B$ remains within an interval of size $\vartheta/2$ for $t \geq t'$. Define $i_{\max} := \arg \max_i \{u_i^{(t')} : 1 \leq i \leq B\}$ (taking the lowest index if this is satisfied by more than one index). Now, for all $1 \leq i' \leq B$ such that $i' \neq i_{\max}$, and for all $t > t'$, we have:

$$u_{i'}^{(t)} - u_{i_{\max}}^{(t)} \leq u_{i'}^{(t')} + \frac{\vartheta}{2} - \left(u_{i_{\max}}^{(t')} - \frac{\vartheta}{2} \right) \leq \vartheta \quad (11)$$

which implies that ρ_1 will not change for $t \geq t'$ (recalling that ϑ is the threshold which must be exceeded before the PASA algorithm will update ρ).

We now proceed via an induction argument. We claim that, provided for all τ_i, h_i and ε_i there exists $\{\eta_j\}_{j=1}^{B+i-1}, H_i$ and T_i such that $\{\rho_j\}_{j=1}^i$ and I_i are ε_i -fixed over τ_i after T_i , then for all τ_{i+1}, h_{i+1} and ε_{i+1} there exists $\{\eta_j\}_{j=1}^{B+i}, H_{i+1}$ and

T_{i+1} such that $\{\rho_j\}_{j=1}^{i+1}$ and I_{i+1} are ε_{i+1} -fixed over τ_{i+1} after T_{i+1} . The case for $i = 1$ we have already shown.

To see our claim holds, suppose we are given values for τ_{i+1} , h_{i+1} and ε_{i+1} . First note that for any $\varepsilon' > 0$, $h' > 0$ we can find η_{B+i} and T' such that, assuming $\{\rho_j\}_{j=1}^i$ remains fixed for all t , \bar{u}_{B+i} will remain in an interval H' of size h' over τ_{i+1} iterations with probability at least $1 - \varepsilon'$ after T' iterations have elapsed (using the same argument as above regarding stochastic approximation algorithms with fixed step-sizes). By assumption, for any τ_i , h_i and ε_i , we can find $\{\eta_j\}_{j=1}^{B+i-1}$ and T_i so that $\{\rho_j\}_{j=1}^i$ will remain the same, and $\sum_{j=1}^{B+i-1} \bar{u}_i$ will remain in an interval of size h_i , for τ_i iterations with probability at least $1 - \varepsilon_i$ in both cases, provided T_i iterations have elapsed. Accordingly we choose h_i and h' such that $h' + h_i \leq \min\{h_{i+1}, \vartheta/2\}$, and we choose τ_i such that $\tau_i \geq T' + \tau_{i+1}$. Hence ρ_{i+1} will remain the same, and $\sum_{j=1}^{B+i} \bar{u}_i$ will remain in an interval of size h_{i+1} , over τ_{i+1} iterations both with probability at least $(1 - \varepsilon_i)(1 - \varepsilon')$ after $T_{i+1} \geq T' + T_i$ iterations have elapsed. For any ε_{i+1} we can choose ε_i and ε' so that $(1 - \varepsilon_{i+1}) > (1 - \varepsilon_i)(1 - \varepsilon')$, and so our claim holds. Hence we can choose η and T such that the vector ρ becomes ε -fixed over τ after T , and so the same holds for $F^{(t)}$. \square

We have taken care to allow the vector η to remain fixed as a function of t . In practical applications, fixed step-sizes will allow an agent to continue to adapt in response to, for example, changes in the environment, and we use fixed step-sizes in our experiments below. Whilst in our experiments in Section 4 we use only a single step-size parameter (as opposed to a vector), the details of the proof point to why there may be merit in using a vector of step-size parameters as part of a more sophisticated implementation of the ideas underlying PASA (i.e. allowing η_k to take on larger values for larger values of the index k , for $k > B$, may allow the algorithm to converge more rapidly).

3 Theoretical analysis

3.1 Uniformly distributed transition function priors

We now set out our main theoretical results. The key idea is that, in many important circumstances, when following a fixed policy an agent will have a tendency to spend nearly all of its time in only a small subset of the state space. We can use this property to our advantage. It means that by focussing on this small area (which is what PASA does) we can eliminate most of the terms which significantly contribute to the expected squared VF error. The trick will be to quantify this tendency. The fact that we have adopted a cell splitting approach will be of critical importance, because it easily permits us to create cells which contain only a single state (allowing us to estimate the VF restricted to single state-action pairs with complete accuracy).

We must make the following assumptions: (1) P is “close to” deterministic, i.e. P , interpreted as an operator with three arguments, can be expressed as follows: $P = (1 - \delta)P_1 + \delta P_2$, where P_1 is a deterministic transition function and P_2 is an arbitrary transition function, and where δ is small (what constitutes “small” will be made clearer below), (2) P has a uniform *prior* distribution, in the sense that, according to our prior distribution for P , the random vector $P(\cdot|s_i, a_j)$ is

independently distributed for all (i, j) and each random variable $P(s_{i'}|s_i, a_j)$ is identically distributed for all (i, j, i') , and (3) π is also “close to” deterministic (i.e. the probability of taking an *off-policy* action is no greater than δ , where an *on-policy* action is the action with the highest probability for a given state, and all other actions are off-policy).

Momentarily putting aside these assumptions, we can make the following observation. If π and P are deterministic, and we pick a starting state s_1 , then the agent will create a path through the state space and will eventually revisit a previously visited state, and will then enter a cycle. Call the set of states in this path (including the cycle) \mathcal{L}_1 and call the set of states in the cycle \mathcal{C}_1 . Denote as L_1 and C_1 the number of states in the path (including the cycle) and the cycle respectively. Of course $L_1 \geq C_1 \geq 1$.

If we now place the agent in a state s_2 (arbitrarily chosen) it will either create a new cycle or it will terminate on the path or cycle created from s_1 . Call \mathcal{L}_2 and \mathcal{C}_2 the states in the second path and cycle (and L_2 and C_2 the respective numbers of states, noting that $C_2 = 0$ is possible if the new path terminates on \mathcal{L}_1 , and in fact that $L_2 = C_2 = 0$ is also possible, if $s_2 \in \mathcal{L}_1$). If we continue in this manner we will have S sets $\{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_S\}$. Call \mathcal{C} the union of these sets, denote as C the number of states in \mathcal{C} , and call C_N the number of sets \mathcal{C}_i which are not empty. We denote as T_i the event that the i th path created in such a manner terminates on itself, and note that, if this does not occur, then $C_i = 0$. In the discussion which follows, until stated otherwise, we assume that (2) holds.

Lemma 1 $E(C_1) = \sqrt{\pi S/8} + O(1)$ and $\text{Var}(C_1) = (32 - 8\pi)S/24 + O(\sqrt{S})$.

Proof Choose any state s_1 . We must have:

$$P(C_1 = i, L_1 = j) = \frac{S-1}{S} \frac{S-2}{S} \dots \frac{S-j+1}{S} \frac{1}{S} = \frac{(S-1)!}{S^j (S-j)!} \quad (12)$$

This means that, for large S , the expected value of C_1 (over the prior distribution for P) can be approximately expressed, making use of Stirling’s approximation $n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$, as:

$$\begin{aligned} E(C_1) &= \sum_{j=1}^S j \sum_{k=j}^S \frac{(S-1)!}{S^k (S-k)!} = \sum_{j=1}^S \frac{j(j+1)}{2} \frac{(S-1)!}{S^j (S-j)!} \\ &= \frac{(S-1)!}{2} \sum_{j=1}^S \frac{(S-j+1)(S-j+2)}{S^{S-j+1} (j-1)!} \\ &= \frac{S!}{2S^{S+1}} \left(\sum_{j=0}^{S-1} \frac{(S-j)^2 S^j}{j!} + \sum_{j=0}^{S-1} \frac{(S-j) S^j}{j!} \right) \\ &= \frac{\sqrt{2\pi S} \left(\frac{S}{e}\right)^S}{2S^{S+1}} \left(\frac{S e^S}{2} + O(\sqrt{S} e^S) \right) = \sqrt{\frac{\pi}{8}} S + O(1) \end{aligned} \quad (13)$$

We have also used the fact that a Poisson distribution with parameter S will, as S becomes sufficiently large, be well approximated by a normal distribution with mean S and standard deviation \sqrt{S} . Hence we can replace the first and second sum in the third equality by the second and first raw moment respectively of a

half normal distribution with variance S . (The error associated with the Stirling approximation is less than order 1.)

Similarly for the variance, we first calculate the expectation of C_1^2 :

$$\begin{aligned} \mathbb{E}(C_1^2) &= \sum_{j=1}^S j^2 \sum_{k=j}^S \frac{(S-1)!}{S^k(S-k)!} = \sum_{j=1}^S \left(\frac{j^3}{3} + \frac{j^2}{2} + \frac{j}{6} \right) \frac{(S-1)!}{S^j(S-j)!} \\ &= \frac{S!}{S^{S+1}} \sum_{j=0}^{S-1} \left(\frac{(S-j)^3}{3} + \frac{(S-j)^2}{2} + \frac{S-j}{6} \right) \frac{S^j}{j!} \\ &= \frac{\sqrt{2\pi S} \left(\frac{S}{e}\right)^S}{S^{S+1}} \left(\sqrt{\frac{2}{\pi}} \frac{2S^{\frac{3}{2}}}{3} e^S + O\left(S e^S\right) \right) = \frac{4}{3}S + O(\sqrt{S}) \end{aligned} \quad (14)$$

As a result:

$$\text{Var}(C_1) = \frac{4}{3}S - \frac{\pi S}{8} + O(\sqrt{S}) = \left(\frac{32 - 8\pi}{24} \right) S + O(\sqrt{S}) \quad (15)$$

□

Note that the expectation can also be derived from the solution to the ‘‘birthday problem’’:⁴ the solution to the birthday problem gives the expectation of L_1 , and since each cycle length (less than or equal to L_1) has equal probability when conditioned on this total path length, we can divide the average by 2.

Lemma 2 $\mathbb{E}(C) < \mathbb{E}(C_1)(\ln S + 1)$ and $\text{Var}(C) \leq O(S \ln S)$.

Proof We will have:

$$\begin{aligned} \mathbb{E}(C) &= \sum_{i=1}^S \mathbb{E}(C_i) = \sum_{i=1}^S \Pr(T_i) \sum_{j=1}^S j \Pr(C_i = j | T_i) \\ &\leq \sum_{i=1}^S \frac{1}{i} \sum_{j=1}^S j \Pr(C_1 = j) < \mathbb{E}(C_1)(\ln S + 1) \end{aligned} \quad (16)$$

And for the variance:

$$\begin{aligned} \text{Var}(C) &= \sum_{i=1}^S \text{Var}(C_i) + 2 \sum_{i=2}^S \sum_{j=1}^{i-1} \text{Cov}(C_i C_j) \leq \sum_{i=1}^S \text{Var}(C_i) \\ &\leq \sum_{i=1}^S \mathbb{E}(C_i^2) = \sum_{i=1}^S \Pr(T_i) \sum_{j=1}^S j^2 \Pr(C_i = j | T_i) \\ &\leq \sum_{i=1}^S \frac{1}{i} \sum_{j=1}^S j^2 \Pr(C_1 = j) < \mathbb{E}(C_1^2)(\ln S + 1) \\ &= \left(\text{Var}(C_1) + \mathbb{E}(C_1)^2 \right) (\ln S + 1) \end{aligned} \quad (17)$$

where we have used the fact that the covariance term must be negative for any pair of lengths C_i and C_j , since if C_i is greater than its mean the expected length of C_j must decrease, and vice versa. □

⁴ For a description of the problem and a formal proof see, for example, page 114 of Flajolet and Sedgewick (2009).

Supposing that π and P are no longer deterministic then we can still define the sets \mathcal{L}_i and \mathcal{C}_i for π and P by considering the *most probable* action and the *most probable* transition. Indeed, if (1) and (3) hold, we can set δ sufficiently low so that the agent will spend an arbitrarily large proportion of its time in \mathcal{C} . Note that if P is deterministic, the transition matrix generated by π and P is not guaranteed to be irreducible or aperiodic, in which case ψ may not exist. This is just a technicality, as a state distribution which is periodic but still restricted to a small number of states does not violate any of our conclusions.

The PASA algorithm, provided X is large enough and provided a subset of the state space has sufficiently high probability and is sufficiently small, will be such that the majority of states with high probability will be represented individually in the final set of basis functions. We can now use this fact, in conjunction with the results we've generated above regarding the distribution of ψ , to demonstrate that PASA will generate a set of basis functions which will be such that the resulting estimate \hat{Q} will have an arbitrarily low error L .

Theorem 1 *For all $\varepsilon_1 > 0$ and $\varepsilon_2 > 0$, there is sufficiently large S and sufficiently small δ such that PASA in conjunction with a suitable RL algorithm will – provided $X \geq K\sqrt{\pi S/8} \ln S \log_2 S$ for some $K > 1$ – generate, with probability no less than $1 - \varepsilon_1$, a VF estimate with $L \leq \varepsilon_2$.*

Proof Using Chebyshev's inequality, and Lemmas 1 and 2, for any $K > 1$ we can choose S sufficiently high so that $C > K\sqrt{\pi S/8} \ln S$ with probability no greater than ε_1 . To see this, take $Y = \sqrt{\pi S/8} \ln S$, $\mu = E(C)$ and $\sigma = \sqrt{\text{Var}(C)}$. We have:

$$\begin{aligned} \Pr(C > KY) &= \Pr(|C - \mu| > KY - \mu) \leq \Pr(|C - \mu| > (K - 1)Y) \\ &\leq \Pr\left(|C - \mu| > (K - 1)\sqrt{\ln S}\sigma\right) \leq \frac{1}{(K - 1)^2 \ln S} \end{aligned} \quad (18)$$

where in the first and second inequalities we assume S is sufficiently large so that all but the highest order terms in Lemma 2 can be ignored.

Since R is bounded and $\gamma < 1$ then for any $\varepsilon > 0$ and S we can also choose δ so that L summed only over states not in \mathcal{C} is no greater than ε_2 . We choose δ so that this is satisfied, but also so that $\psi_i > \sum_{i': s_{i'} \notin \mathcal{C}} \psi_{i'}$ for all elements of $\{\psi_i : s_i \in \mathcal{C}\}$. This is possible since each such ψ_i will be bounded from below for all $\delta > 0$. Now provided that $C \log_2 S \leq X$ then each state in \mathcal{C} will eventually be in its own cell. The RL algorithm will have no error for each such state so L will be no greater than ε_2 . \square

We can prove in almost identical fashion that, under the same assumptions, both (i) MSE and (ii) L redefined so that all actions are weighted equally, i.e. equation (3) with the factor π removed, will similarly be arbitrarily close to zero with arbitrarily high probability.

The bound on X provided represents a significant reduction in complexity when S starts to take on a size comparable to many real world problems, and could make the difference between a problem being tractable and intractable. It also seems likely that the bound on X in Theorem 1 can be improved upon (see the discussion in Section 4). Conditions (1) and (3) are commonly encountered in practice, in particular (3) which can be taken to reflect a “greedy” policy. Condition

(2) can be interpreted as the transition function being “completely unknown” (and can be generalised somewhat, see below).

We might ask what happens when π is no longer generated randomly, but rather, for example, according to a process of policy iteration. It is intuitively plausible that this should typically have the effect of reducing $E(C)$. We will not offer a general proof. However note that, if for example $\zeta = \alpha S$ states have reward of one, for $\alpha \in (0, 1)$, and all other states have reward of zero, the agent must determine a shorter path on average than a random path between each pair of such states (otherwise its performance will be no better than a random policy). Denote as d the average path length between reward states determined by an agent which is updating its policy so as to increase reward. By the reasoning we’ve just stated d will not increase above a bounded value as S increases. We can then apply the same reasoning as the random policy case to conclude that at most $O(\sqrt{\zeta} \ln \zeta)$ such states will be visited under such a policy on average, and that $E(C_1) \leq O(d\sqrt{\zeta} \ln \zeta) \leq O(\sqrt{S} \ln S)$.

On a final related note, *fixed* state aggregation will *also* tend to have zero error if $X > O(\sqrt{S} \ln S)$ and all our other assumptions continue to hold. This is because the probability of more than one state falling into a cell tends to zero if the number of cells grows at a faster rate than the size of \mathcal{C} . Of course the probability of all states falling into a unique cell increases much more slowly than is true for PASA as S is increased. Alterations can be made to the PASA algorithm (with no impact on computational complexity) that can remove the $\log_2 S$ factor in Theorem 1. However such an alternative algorithm is more complex to describe and unlikely to perform noticeably better in a practical setting.

3.2 Generalisations to other prior distributions

There is some scope to extend our results beyond uniform transition function priors. Some knowledge about the value of P (i.e. a non-uniform prior) can have important implications for our results. If we know, for example, that $P(s_{i+1}|s_i, a_j) = 1$ for all $1 \leq i \leq S - 1$ and all j , and $P(s_1|s_S, a_j) = 1$ for all j , then $C_1 = S$, and Lemma 1 clearly doesn’t hold. There are, however, more general sets of priors for which the results will hold.

Our interest is in how the prior influences the values of the moments which we calculated (for uniform priors) in Lemmas 1 and 2. We again assume a deterministic transition function, and a fixed deterministic policy π . In our discussion around uniform priors, we were also able to assume that the sequence of states s_1, \dots, s_S used to generate the sets $\mathcal{C}_1, \dots, \mathcal{C}_S$ were selected arbitrarily. In a more general setting we may need to assume that this sequence is generated according to some specific probability distribution. Since, in the discussion below, we will generally assume that the prior distribution of the transition probabilities is identical for different actions, we will be able to continue to assume that π is arbitrary.⁵

Appealing to techniques which make use of the notion of Schur convexity (Marshall et al 2011), it’s possible to show that, if the random vector $P(\cdot|s_i, a_j)$ is independently distributed for all (i, j) , and $\Pr(P(s_{i''}|s_i, a_j) = 1) = \Pr(P(s_{i''}|s_{i'}, a_{j'}) =$

⁵ In the generalisations we discuss, we could potentially remove some assumptions from P at the expense of imposing heavier constraints on π .

1) for all (i, j, i', j', i'') , then, given an arbitrary policy π , and assuming the sequence of starting states s_1, \dots, s_S are distributed uniformly at random, $E(L_1)$ and $\text{Var}(L_1)$ are minimised where the prior for P is uniform. Using this fact, Theorem 1 can be extended to such priors (we omit the details, though the proof uses arguments substantially equivalent to those used in Theorem 1). If we continue to assume an arbitrary policy and that the starting states are selected uniformly at random, we can consider a yet more general class of priors, using a result from Karlin and Rinott (1984). Their result can be used to demonstrate that, of the set of priors which satisfy the following three conditions – (i) that $\Pr(P(s_{i'}|s_i, a_j) = 1) = \Pr(P(s_{i'}|s_i, a_{j'}) = 1)$ for all (i, j, i', j') , (ii) that:

$$\begin{aligned} \Pr(P(s_{i_3}|s_{i_1}, a_j) = 1) &> \Pr(P(s_{i_4}|s_{i_1}, a_j) = 1) \\ &\Rightarrow \Pr(P(s_{i_3}|s_{i_2}, a_j) = 1) > \Pr(P(s_{i_4}|s_{i_2}, a_j) = 1) \end{aligned} \quad (19)$$

for all (i_1, i_2, i_3, i_4, j) , and (iii) the random vector $P(\cdot|s_i, a_j)$ is independently distributed for all (i, j) – the uniform prior will again maximise the values $E(L_1)$ and $\text{Var}(L_1)$. Since $C_i \leq L_i$, an equivalent result to Theorem 1 (though not necessarily with the same constant $\sqrt{\pi/8}$) can similarly be obtained for this even larger set of priors (we again omit the details and note that the arguments are substantially equivalent to those in Theorem 1).

Both these results assume a degree of similarity in the transition prior probabilities for each state. A perhaps more interesting potential generalisation is as follows. We can define a *balanced* prior for a deterministic transition function P as any prior such that the random vector $P(\cdot|s_i, a_j)$ is independently distributed for all (i, j) , and we have $\Pr(P(s_{i'}|s_i, a_j) = 1) = \Pr(P(s_{i'}|s_i, a_{j'}) = 1)$, $\Pr(P(s_{i'}|s_i, a_j) = 1) = \Pr(P(s_i|s_{i'}, a_j) = 1)$ and $\Pr(P(s_i|s_i, a_j) = 1) \geq \frac{1}{S}$ for all (i, j, i', j') . In essence, the prior probability of transitioning from state s_i to $s_{i'}$ is the same as transitioning in the reverse direction from $s_{i'}$ to s_i . This sort of prior would be reflective of many real world problems which incorporate some notion of a geometric space with distances, such as navigating around a grid – and as such represents an important generalisation. The difference to the uniform prior is that we now have an expectation that some states will be “close” to one-another, and other states will be further apart. However similar to the uniform prior there is no inherent “flow” creating cycles which have larger expected value than C in the uniform case.

It is not hard to conceive of examples where $E(C_1)$ may be significantly reduced for a particular balanced prior compared to the uniform case. It would furthermore appear plausible that, amongst the set of all balanced priors, $E(L_1)$ would be maximised for the uniform prior. Indeed investigation using numerical optimisation techniques demonstrates this is the case for $S \leq 8$, even when a fixed arbitrary starting state is selected relative to the balanced set of transition probabilities. The techniques used for the generalisations stated above cannot be used to prove a similar result for balanced priors.⁶ We conjecture, based on our numerical analysis,

⁶ The earlier stated results follow in both cases from the stronger statement that $\Pr(L_1 > k)$ is maximised for all k by a uniform prior, from which our conclusions regarding the moments follow. In the case of balanced priors such a strong result does not hold, which can be seen by, for large S , and taking $k = S$, comparing a uniform prior to a prior where all transitions outside of a single fixed cycle covering all S states have probability zero, and where the prior probability of a transition in either direction along this cycle is $(1 - 1/S)/2$.

that the uniform prior does maximise $E(L_1)$ for all S , which would carry the implication, since $C_1 \leq L_1$, that Lemma 1 can be used to argue $E(C_1) \leq O(\sqrt{S})$ and $\text{Var}(C_1) \leq O(S)$.

Even if this conjecture holds, we cannot extend Theorem 1 to balanced priors, which we can see with a simple example. Set $\Pr(P(s_i|s_i, a_j) = 1) \geq 1 - \varepsilon$ for all (i, j) where ε is small. Provided that the transition matrix associated with π and P is irreducible, then $C = S$ for all S .

Notwithstanding that a formal result equivalent to Theorem 1 is unavailable for balanced priors, by selecting suitable parameter values we should still be able to exploit the apparent tendency of the agent to spend a majority of the time in a small subset of the state space. The main difference is that this small subset may change slowly over time. The situation may become further altered in our favour once policy iteration is introduced, and policies are no longer random but rather target states with reward. Our experimental results focus on transition functions drawn from a uniform distribution, however further research could help indicate the extent to which PASA can generate a VF estimate with low error, subject to the slightly altered dynamics introduced by balanced priors.

4 Experimental results

Our main objectives in this section are to: (a) Test empirically the tendency for an agent to spend the majority of its time in a small subset of the state space, and (b) Conduct an experimental comparison of PASA to fixed state aggregation, to demonstrate that, on average, PASA will help generate VF estimates with lower MSE (and therefore that the theoretical advantages to unsupervised basis function adaptation can be realised in a practical setting).

4.1 Empirical results relating to the stationary state distribution

Whilst the bound we placed on X in Theorem 1 may be of help when S is very large, note that if we set $X = \sqrt{\pi S/8} \ln S \log_2 S$ (i.e. such that we have a fifty per cent chance of being guaranteed to be able to represent every high probability state), then we will have $X > S$ up until the point where $S \geq 3,748$. It appears likely that the bound in Theorem 1 is a loose bound. We conduct some simple experiments to examine how C behaves, in particular for lower values of S .

In Figure 2 we report on the outcome of a sequence of 1,000 independent trials where C was calculated explicitly, in order to test our theoretical bounds. The results demonstrate that the estimates for $E(C_1)$ and $\text{Var}(C_1)$ are accurate, but that the bound in Lemma 2 is (perhaps not surprisingly) generous. In Figure 2 we have also run equivalent tests on a square grid world environment⁷ (we have used C_1^g , C_N^g and C^g to distinguish these values from those obtained for the uniform prior). This helps to reinforce the discussion around balanced transition function priors (we can see that C_1^g is very small compared to the uniform case, and may

⁷ In our experiment we actually set $\Pr(s_i|s_i, a_j) = 0$ for all i and j (the agent is forced to move in one of four directions: up, down, left or right). This violates the exact balanced definition, however would only serve to increase $E(C_1^g)$.

in fact be bounded as a function of S , however that C_N^g , and therefore the value C^g , appears to increase linearly as a function of S).

We can also examine the empirical distribution of ψ given a sample trajectory, to see the extent to which our Section 3 assumptions do indeed result in only a few states having high probability. This is worth exploring since the value of δ strictly required by Theorem 1 may be too low to be realistic in practice (i.e., given a too-large value for δ , we might expect to see the stationary distribution “leak” to include states outside \mathcal{C}). Assume that the agent follows an ϵ -greedy policy. We define:

$$M_p = \min \left\{ |\mathcal{R}| : \mathcal{R} \subseteq \{1, 2, \dots, S\}, \sum_{i \in \mathcal{R}} \psi_i \geq p \right\} \quad (20)$$

Hence M_p is the minimum number of states required to amass a proportion $p \in [0, 1]$ of the total probability of the stationary state distribution (as p approaches 1 and δ approaches 0, the value can be considered as a rough experimental equivalent of C). Figure 3 demonstrates that, when ϵ is held constant as S is increased, we start to see what appears to be a linear increase in M_p as a function of S . It seems likely that this is because of leakage into states outside \mathcal{C} . As C/S becomes smaller and ϵ remains the same, this leakage will become more pronounced. Whilst this doesn’t invalidate any of the underlying principles, in a practical setting it is something to remain aware of. Even where ϵ is reasonably high, e.g. 0.01, and where we require a high value of p , e.g. $p > 0.99$, only a reasonably small number of states is required to obtain a total probability of p , e.g. around 20% of states for $S = 1,000$, implying that in general we can set X comparatively low.

The implication of these results is that we expect we should be able to use PASA to effectively reduce error in the VF estimate in a practical setting (including where S is low, for example less than 1,000). This is what we next investigate. As a final comment, whilst we need $X \geq C \log_2(S)$ to theoretically guarantee capturing every high-probability state individually, the factor of $\log_2(S)$ is not likely to be essential in a practical setting (as it covers the worst case scenario in terms of how the high-probability states are arranged).

4.2 Comparison of PASA to fixed state aggregation

We have run a series of experiments to test the performance of PASA compared with simple fixed state aggregation. Before stating the results it will be helpful to make some comments first.

If γ approaches 1 and δ is too large as a function of γ , then we will find that all values $Q^\pi(s_i, a_j)$ will approach the same value, for any reward function R . This is because the discounted reward is heavily weighted into the future, and due to randomness in the agent’s trajectory, future reward will not be heavily impacted by each individual decision. Hence, under such conditions, both MSE and L for fixed state aggregation will *also* tend towards arbitrarily low values. Furthermore, if, given a state aggregation architecture, the number of high probability states in each cell is large, then the average value for each state-action over each cell will also approach the same value (this is due to the law of large numbers). The second point is important when calculating L , since if F is such that all cells are large (and particularly if γ is close to 1) then each estimate $\hat{Q}(s_i, a_j)$ will tend to

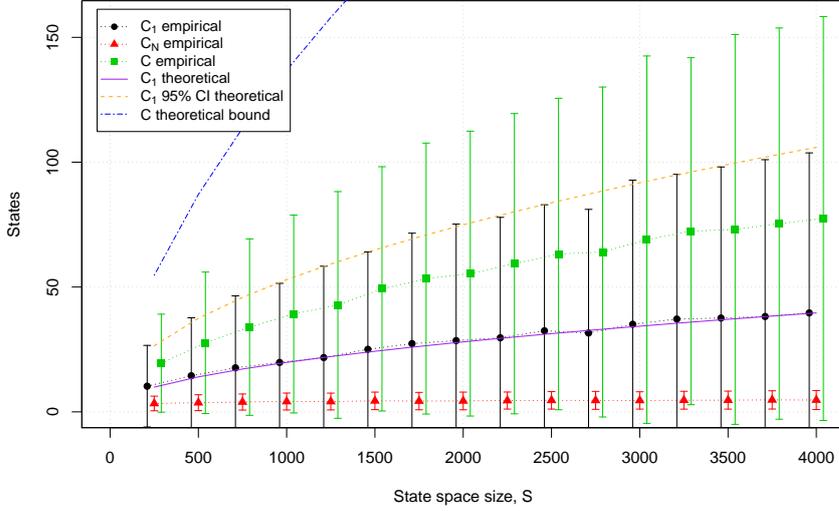
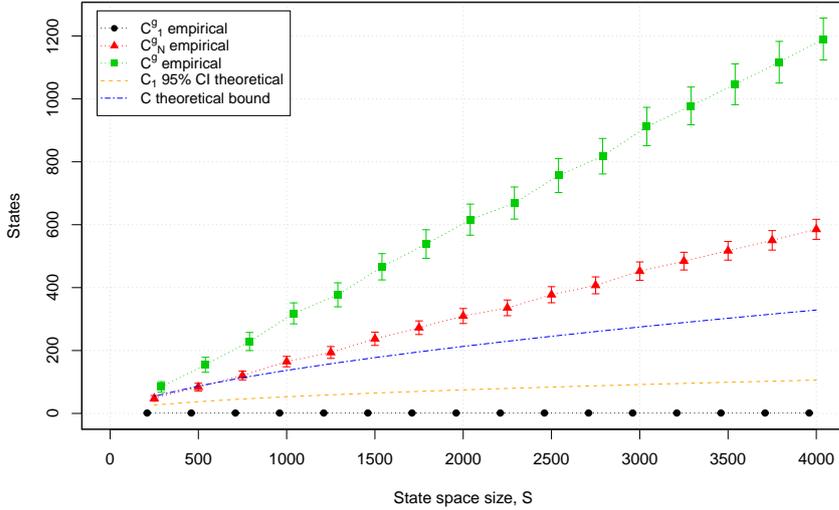
(a) Uniformly distributed P (b) Grid-world P

Fig. 2: The agent spends the majority of its time in a small number of states. In this chart we have explicitly calculated the *average* values C , C_1 and C_N for 1,000 independent randomly generated transition functions. Chart (a) is for the uniform transition function prior. We have included error bars which reflect 95% confidence intervals (the confidence intervals assume the data points are normally distributed in each case, which is true for C_N and C , though less so for C_1). Chart (b) is for a grid-world environment (the theoretical bounds for the uniform case are reproduced for reference).

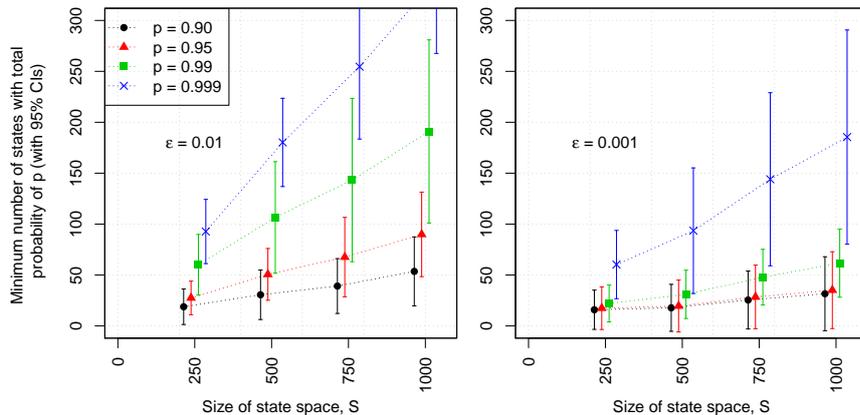


Fig. 3: The agent spends the majority of its time in a small number of states. The states are ordered in terms of their stationary probability ψ_i , and we then determine the average (over 100 trials) of M_p for $p = 0.9$, $p = 0.95$, $p = 0.99$ and $p = 0.999$. Confidence intervals (95%) are included.

a single mean value, and L will take on a low value (even if MSE is high, and the VF estimate is poor).

For these reasons, we have taken care not to set γ too close to 1 (keeping in mind the value of δ), and we have used MSE to compare PASA to fixed state aggregation. Since MSE is a more accurate measure, this will only serve to strengthen the results. The only downside is that, since an exact VF estimate is impossible to calculate for large S , we can only perform our comparison for relatively low values of S . We then use our theoretical insights to infer that similar performance differences will exist for large S . In all cases we have reported the square root of the MSE (as this provides a less distorted comparison).

The two algorithms we will be comparing are PASA combined with a SARSA⁸ RL algorithm, and a SARSA RL algorithm with fixed state aggregation, where the cells are sized as equally as possible. It can be shown that SARSA will generate a VF estimate with minimal MSE given a fixed policy and a fixed linear approximation architecture (Bertsekas and Tsitsiklis 1996).

We have tested both algorithms with sequences of randomly generated environments (where P is deterministic and with a uniform prior, the prior for R is such that the reward for each state-action pair has an independent standard normal distribution, and π is ϵ -greedy), and then estimated (by sampling) the MSE of the estimate VF generated by the RL algorithm.

⁸ Further details regarding the SARSA algorithm can be found, for example, in Sutton and Barto (1998). For the SARSA component we have used a constant step-size of 3×10^{-4} which we then weighted by the inverse of the probability of the action selected. Such normalisation is not strictly necessary, however it allows the VF estimate for infrequently selected actions to converge more quickly. The random policy π was chosen such that the same action a_j is the “preferred” action (the action taken with probability $1 - \epsilon$) for all states in each particular cell.

Table 1: Summary of results of PASA and fixed state aggregation MSE.

Experiment	S	X	$\sqrt{\text{MSE}}$ (fixed) ^a		$\sqrt{\text{MSE}}$ (PASA)		% $\sqrt{\text{MSE}}$ decrease	
			$\gamma = 0$	$\gamma = 0.99$	$\gamma = 0$	$\gamma = 0.99$	$\gamma = 0$	$\gamma = 0.99$
1	250	70	0.19	7.5	0.06	3.23	67.7%	56.9%
2	500	100	0.26	9.97	0.09	4.42	63.2%	55.6%
3	750	130	0.23	7.99	0.09	4.33	59.4%	45.8%
4	1,000	160	0.27	9.41	0.12	5.42	55.9%	42.4%

^a We take the average MSE over the last fifth of each trial.

For each algorithm we ran eight separate experiments, for S set at 250, 500, 750 and 1,000, and for $\gamma = 0.99$ and $\gamma = 0$, for 100 randomly generated trials each.⁹ In all cases both algorithms were tested against an identical set of environments. For each value of S listed above we set X equal to 70, 100, 130 and 160 respectively (and B equal to 35, 50, 65 and 80), guided in part by Figure 3 and the other comments above (but also in part heuristically).¹⁰ The results of these experiments are shown in Table 1 (and Figure 4 for $\gamma = 0.99$). The values of other parameters used were as follows: $A = 2$, $\epsilon = 0.001$, $\vartheta = 0.001$, $\nu = 10,000$ and $\eta = 1 \times 10^8$.

With respect to sampling the MSE, we compared the VF estimate \hat{Q} for each state visited to the exact Q^π for the policy π , which for each experiment was calculated in advance exactly by solving the system of equations involving θ using least squares. Quite a large number of iterations were used in each trial (500 million). This was in part to accommodate the most complex environment ($S = 1,000$), and also to provide clear evidence that the estimates have stabilised. Such long trials are not necessary to see a significant difference between the two algorithms (see Figure 4). Minimising the number of iterations required to see a strong difference between PASA and fixed state aggregation would depend on optimising parameters like ϑ , ϵ and η subject to the choices for γ and S .

4.3 Comments on experimental results

Since the theoretic results assume, for example, that η and ϵ can be made arbitrarily low, in a practical setting we won't see MSE falling to zero. However Table 1 demonstrates that, at a low computational cost, we are able to significantly reduce MSE from the fixed state aggregation case. The table shows us that for $\gamma = 0.99$ error can be reduced by around 40 to 60 per cent, and for the less complex case of $\gamma = 0$ by as much as approximately 50 to 70 per cent. As noted above, the effect of PASA becomes slightly less pronounced as γ approaches 1 (in particular if ϵ is comparatively large). Figure 4 shows that, as S increases, PASA can make the algorithm take slightly longer to generate a VF with low MSE (which might

⁹ For $S = 1,000$ the running time was slightly less than seven hours for a single trial on an Intel(R) Xeon(R) CPU E5-4650 0 @ 2.70GHz for both algorithm variants, although this includes the time to calculate beforehand an exact solution (to allow for MSE to be calculated).

¹⁰ Setting B to approximately half X appears to be a reliable way of ensuring that a large number of high probability states fall into a singleton cell. Smaller values of B tend to exhaust many cells in seeking to isolate (into smaller and smaller cells) one or two individual states, whereas greater values of B may mean that PASA has too little control over the final set of basis functions.

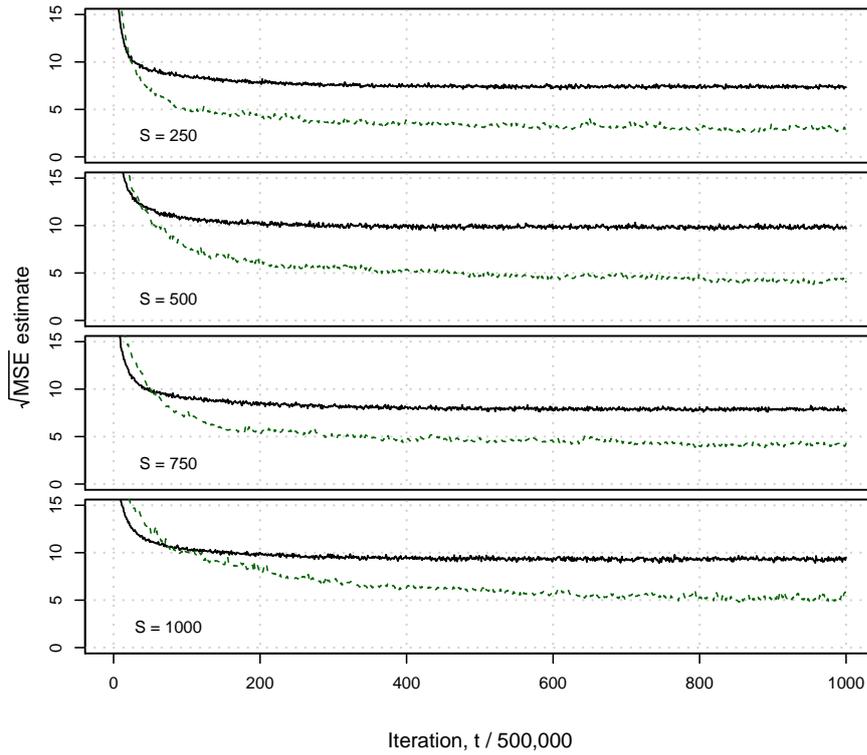


Fig. 4: Comparative MSE for PASA (green) versus fixed state aggregation (black) as a function of t (for $\gamma = 0.99$). Each line is an average over 100 independent trials. MSE is averaged over blocks of 500 thousand iterations. As S becomes larger, PASA can tend to take longer to decrease MSE, as the mapping F requires some time to converge (this is most pronounced for $S = 1,000$).

be expected since the PASA component needs time to converge), however in the longer term MSE is significantly reduced.

Optimising parameters such as η , ϵ or even the value of X (which we've set well below the bound in Section 3) might be expected to further increase the disparity between PASA and fixed state aggregation (the latter cannot be improved since it has no parameters to adjust outside the underlying RL algorithm). In the context of policy improvement, even small reductions in the VF error may have a large impact on achieving reward.

5 Conclusion

As we noted in the introduction, whilst basis function adaptation methods have been the subject of a large amount of research recently, unsupervised methods have

not received much attention (compared to methods which seek to estimate VF error explicitly). This may be an oversight. As we have seen there are many circumstances reflecting real-world problems where unsupervised methods can be very effective in creating an approximation architecture, and thereby help to generate more accurate VF estimates (both in theory and in practice). What distinguishes unsupervised methods from more complex alternatives, however, is their simplicity. As we have seen, an algorithm such as PASA carries only minimal additional costs (both in terms of sampling and computational demands) compared to the RL algorithm which it supports. It is their effectiveness combined with their low cost which, in our view, make such techniques a promising candidate for further research.

In the setting of policy improvement the advantages offered by unsupervised methods have the potential to be particularly important. Each policy update requires a VF estimate (where the VF estimate is based on the current policy) however each policy update will change the VF (perhaps significantly), requiring a new estimate. So accurately and *efficiently* estimating Q^π for each new policy π is critical. Some initial experimentation suggests that the PASA algorithm can have a significant impact on RL algorithm performance where policy updates are introduced. This should perhaps not come as a surprise given the reduced error in \hat{Q} shown in Section 4.

The PASA algorithm represents only a relatively naive or simplistic application of the idea of unsupervised basis function adaptation (i.e. the idea of using visit frequency to guide which parts of the VF to represent in more detail). It appears likely that a number of improvements could be made to the algorithm in order to further optimise its performance.¹¹

From a theoretical perspective the assumption around scoring functions being weighted by ψ is crucial. The nature of the VF estimate generated by PASA and its associated RL algorithm is that the VF will be well estimated for states which are frequently visited under the existing policy. However this results in poorer estimates of the value of deviating from the current policy. Thus, even though the expected VF error may be low, it does not immediately follow that an algorithm can use this to optimise its policy via standard policy iteration (since the consequences of deviating from the current policy are less clearly represented). Ultimately, however, the theoretical implications of the improved VF estimate in the context of policy iteration are complex, and would need to be the subject of further research.

References

1. Bernstein, A., & Shimkin, N. (2010). Adaptive-resolution reinforcement learning with efficient exploration in deterministic domains. *Machine Learning*, 81(2), 359-397.
2. Bertsekas, D., & Tsitsiklis, J. (1996). *Neuro-Dynamic Programming*. Athena Scientific.
3. Bertsekas, D., & Yu, H. (2009). Convergence results for some temporal difference methods based on least squares. *IEEE Transactions on Automatic Control*, 54(7), 1515-1531.

¹¹ For example, given that the number of high probability states in each initial cell \mathcal{X} is binomially distributed, and this number will have low variance if B is substantially lower than $E(C)$, it seems likely that, in choosing a cell to split, we could initially only consider a subset of the set of all available cells. This would potentially allow PASA to converge more quickly, and be less likely to change as a result of random fluctuation in the value of \bar{u} .

4. Bucklew, J., Kurtz, T., & Sethares, W. (1993). Weak convergence and local stability properties of fixed step size recursive algorithms. *IEEE Transactions on Information Theory*, 30(3), 966-978.
5. Buşoniu, L., Babuška, R., De Schutter, B., & Ernst, D. (2009). *Reinforcement Learning and Dynamic Programming Using Function Approximators*. Taylor & Francis CRC Press.
6. Di Castro, D., & Mannor, S. (2010). Adaptive bases for reinforcement learning. *49th IEEE Conference on Decision and Control*, 312-327.
7. Flajolet, P., & Sedgewick, R. (2009). *Analytic Combinatorics*. Cambridge University Press.
8. Karlin, S., & Y. Rinott (1984). Random replacement schemes and multivariate majorization. *Lecture Notes-Monograph Series*, 35-40.
9. Mahadevan, S., Giguere, S., & Jacek, N. (2013). Basis adaptation for sparse nonlinear reinforcement learning. *Proceedings of the Conference on Artificial Intelligence*.
10. Marshall, A. W., Olkin, I., & Arnold, B. C. (2011). *Inequalities: Theory of Majorization and Its Applications*. Springer Series in Statistics.
11. Menache, I., Mannor, S., & Shimkin N. (2005). Basis function adaptation in temporal difference reinforcement learning. *Annals of Operations Research*, 134(1), 215-238.
12. Moore, A., & Atkeson, C. (1995). The parti-game algorithm for variable resolution reinforcement learning in multidimensional state-spaces. *Machine Learning*, 21(3), 199-233.
13. Munos, R., & Moore, A. (2002). Variable resolution discretization in optimal control. *Machine Learning*, 29(2-3), 291-323.
14. Singh, S., Jaakkola, T., & Jordan, M. (1995). Reinforcement learning with soft state aggregation. *Advances in Neural Information Processing Systems*, 361-368.
15. Sutton, R., & Barto, A. (1998). *Reinforcement Learning: An Introduction*. MIT Press.
16. Whiteson, S., Taylor, M., & Stone, P. (2007). Adaptive tile coding for value function approximation. University of Texas at Austin, Technical Report.