

# Learning to Randomize and Remember in Partially-Observed Environments

Radford M. Neal, University of Toronto

Dept. of Statistical Sciences and Dept. of Computer Science

<http://www.cs.utoronto.ca/~radford>

- I. Background on Reinforcement Learning with Fully Observed State
- II. Learning Stochastic Policies When the State is Partially Observed
- III. Learning What to Remember of Past Observations and Actions
- IV. Can This Work For More Complex Problems?

# The Reinforcement Learning Problem

Typical “supervised” and “unsupervised” forms of machine learning are very specialized compared to real-life learning by humans and animals:

- We seldom learn based on a fixed “training set”, but rather based on a continuous stream of information.
- We also act continuously, based on what we’ve learned so far.
- The effects of our actions depend on the state of the world, of which we observe only a small part.
- We obtain a “reward” that depends on the state of the world and our actions, but aren’t told what action would have produced the most reward.
- Our computational resources (such as memory) are limited.

The field of *reinforcement learning* tries to address such realistic learning tasks.

# Formalizing a Simple Version of Reinforcement Learning

Let's envision the world going through a sequence of *states*,  $s_0, s_1, s_2, \dots$ , at integer times. We'll start by assuming that there are a finite number of possible states.

At every time, we take an *action* from some set (assumed finite to begin with).

The sequence of actions taken is  $a_0, a_1, a_2, \dots$

As a consequence of the state,  $s_t$ , and action,  $a_t$ , we receive some *reward* at the next time step, denoted by  $r_{t+1}$ , and the world changes to state  $s_{t+1}$ .

Our aim is to maximize something like the total “discounted” reward we receive over time.

The discount for a reward is  $\gamma^{k-1}$ , where  $k$  is the number of time-steps in the future when it is received, and  $\gamma < 1$ . This is like assuming a non-zero interest rate — money arriving in the future is worth less than money arriving now.

# Stochastic Worlds and Policies

The world may not operate deterministically, and our decisions also may be stochastic. Even if the world is really deterministic, an imprecise model of it will need to be probabilistic.

We assume the *Markov property* — that the future depends on the past only through the present state (really the definition of what the state is).

We can then describe how the world works by a transition/reward distribution, given by the following probabilities (assumed the same for all  $t$ ):

$$P(r_{t+1} = r, s_{t+1} = s' \mid s_t = s, a_t = a)$$

We can describe our own *policy* for taking actions by action probabilities (again, assumed the same for all  $t$ , once we've finished learning a policy):

$$P(a_t = a \mid s_t = s)$$

This assumes that we can observe the entire state, and use it to decide on an action. Later, I will consider policies based on partial observations of the state.

# Exploration Versus Exploitation

If we know exactly how the world works, and can observe the entire state of the world, there is no need to randomize our actions — we can just take an optimal action in each state.

But if we don't have full knowledge of the world, always taking what appears to be the best action might mean we never experience states and/or actions that could produce higher rewards. There's a tradeoff between:

*exploitation*: seeking immediate reward

*exploration*: gaining knowledge that might enable higher future reward

In a full Bayesian approach to this problem, we would still find that there's always an optimal action, accounting for the value of gaining knowledge, but computing it might be infeasible. A practical approach is to *randomize* our actions, sometimes doing apparently sub-optimal things so that we learn more.

# The $Q$ Function

The expected total discounted future reward if we are in state  $s$ , perform an action  $a$ , and then follow policy  $\pi$  thereafter is denoted by  $Q^\pi(s, a)$ .

This  $Q$  function satisfies the following consistency condition:

$$Q^\pi(s, a) = \sum_r \sum_{s'} \sum_{a'} P(r_{t+1} = r, s_{t+1} = s' \mid s_t = s, a_t = a) P^\pi(a_{t+1} = a' \mid s_{t+1} = s') (r + \gamma Q^\pi(s', a'))$$

Here,  $P^\pi(a_{t+1} = a' \mid s_{t+1} = s')$  is an action probability determined by the policy  $\pi$ .

If the optimal policy,  $\pi$ , is deterministic, then in state  $s$  it must clearly take an action,  $a$ , that maximizes  $Q^\pi(s, a)$ .

So knowing  $Q^\pi$  is enough to define the optimal policy. Learning  $Q^\pi$  is therefore a way of learning the optimal policy without having to learn the dynamics of the world — ie, without learning  $P(r_{t+1} = r, s_{t+1} = s' \mid s_t = s, a_t = a)$ .

# Exploration While Learning a Policy

When we don't yet know an optimal policy, we need to trade off between exploiting what we do know versus exploring to obtain useful new knowledge.

One simple scheme is to take what seems to be the best action with probability  $1 - \epsilon$ , and take a random action (chosen uniformly) with probability  $\epsilon$ . A larger value for  $\epsilon$  will increase exploration.

We might instead (or also) randomly choose actions, but with a preference for actions that seem to have higher expected reward — for instance, we could use

$$P(a_t = a \mid s_t = s) \propto \exp(Q(s, a) / T)$$

where  $Q(s, a)$  is our current estimate of the  $Q$  function for a good policy, and  $T$  is some “temperature”. A larger value of  $T$  produces more exploration.



# Learning a $Q$ Function and Policy with 1-Step SARSA

Recall the consistency condition for the  $Q$  function:

$$Q^\pi(s, a) = \sum_r \sum_{s'} \sum_{a'} P(r_{t+1} = r, s_{t+1} = s' \mid s_t = s, a_t = a) P^\pi(a_{t+1} = a' \mid s_{t+1} = s') (r + \gamma Q^\pi(s', a'))$$

This suggests a Monte Carlo approach to incrementally learning  $Q$  for a good policy. At time  $t+1$ , after observing/choosing the states/actions  $s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}$  (hence the name SARSA), we update our estimate of  $Q(s_t, a_t)$  for a good policy by

$$Q(s_t, a_t) \leftarrow (1-\alpha) Q(s_t, a_t) + \alpha (r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}))$$

Here,  $\alpha$  is a “learning rate” that is slightly greater than zero.

We can use the current  $Q$  function and the exploration parameters  $\epsilon$  and  $T$  to define our current policy:

$$P(a_t = a \mid s_t = s) = \frac{\epsilon}{\#\text{actions}} + (1-\epsilon) \frac{\exp(Q(s, a) / T)}{\sum_{a'} \exp(Q(s, a') / T)}$$

- I. Background on Reinforcement Learning with Fully Observed State
- II. Learning Stochastic Policies When the State is Partially Observed
- III. Learning What to Remember of Past Observations and Actions
- IV. Can This Work For More Complex Problems?

# Learning in Environments with Partial Observations

In real problems we seldom observe the full state of the world. Instead, at time  $t$ , we obtain an observation,  $o_t$ , related to the state by an observation distribution,

$$P(o_t = o \mid s_t = s)$$

This changes the reinforcement learning problem fundamentally:

- 1) Remembering past observations and actions can now be helpful.
- 2) If we have no memory, or only limited memory, an optimal policy must sometimes be stochastic.
- 3) A well-defined  $Q$  function exists only if we assume that the world together with our policy is ergodic.
- 4) We cannot in general learn the  $Q$  function with 1-Step SARSA.
- 5) An optimal policy's  $Q$  function is not sufficient to determine what action that policy takes for a given observation.

Points (1) – (3) above have been known for a long time (eg, Singh, Jaakola, and Jordan, 1994). Point (4) seems to have been at least somewhat appreciated.

Point (5) initially seems counter-intuitive, and doesn't seem to be well known.

# Memoryless Policies and Ergodic Worlds

To begin, let's assume that we have no memory of past observations and actions, so a policy,  $\pi$ , is specified by a distribution of actions given the current observation,

$$P^\pi(a_t = a \mid o_t = o)$$

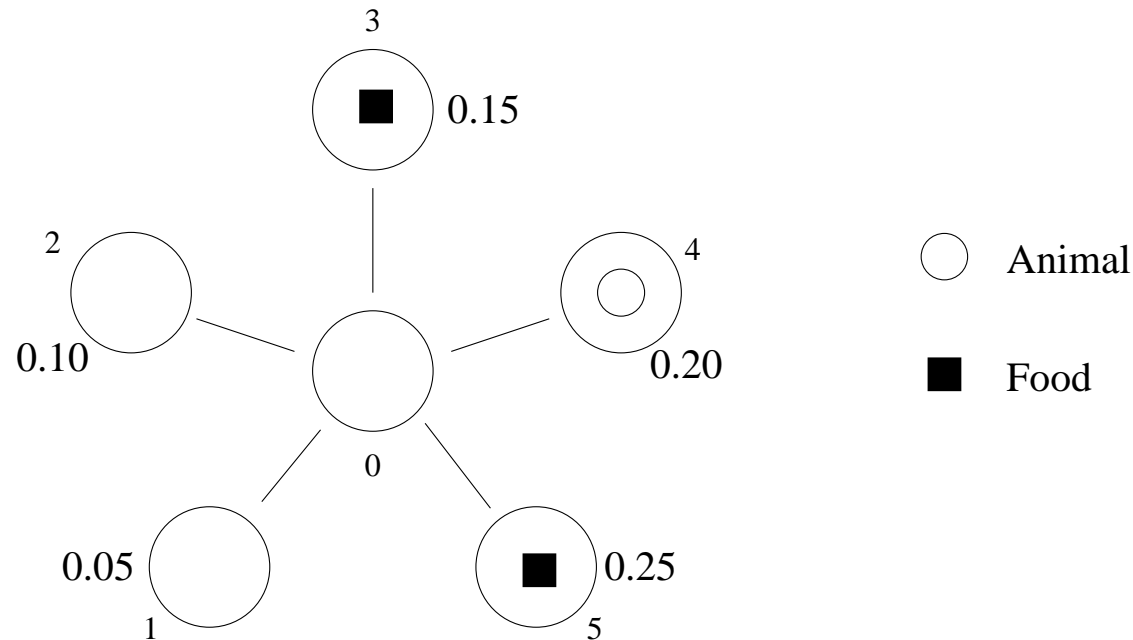
We'll also assume that the world together with our policy is *ergodic* — that all actions and states of the world occur with non-zero probability, starting from any state. In other words, the past is eventually “forgotten”.

This is partly a property of the world — that it not become “trapped” in a subset of the state space, for any sequence of actions we take.

If the world is ergodic, a sufficient condition for our policy is that it give non-zero probability to all actions given any observation. We may want this anyway for exploration.

# Grazing in a Star World: A Problem with Partial Observations

Consider an animal grazing for food in a world with 6 locations, connected in a star configuration:



The centre point (0) never has food. Each time step, food grows at an outer point (1, ..., 5) that doesn't already have food with probabilities shown above. When the animal arrives at a location, it eats any food there. Each time step, it can move along one of the lines shown, or stay where it is.

The animal can observe where it is (one of 0, 1, ..., 5), but not where food is.

Reward is +1 if food is eaten, -1 if attempts invalid move (goes to 0), 0 otherwise.

## Defining a $Q$ Function of Observation and Action

We'd like to define a  $Q$  function using observations rather than states, so that  $Q(o, a)$  is the expected total discounted future reward from taking action  $a$  when we observe  $o$ .

Note! This makes sense only if we assume ergodicity — otherwise  $P(s_t = s \mid o_t = o)$ , and hence  $Q(o, a)$ , are not well-defined.

Also...

- $Q(o, a)$  will depend on the policy followed in the past, since the past policy affects  $P(s_t = s \mid o_t = o)$ .
- $Q(o, a)$  will *not* be the expected total discounted future reward *conditional* on events in the recent past, since the future is not independent of the past given only our current observation (rather than the full state at the current time).
- But with an ergodic world + policy,  $Q(o, a)$  will approximate the expected total discounted future reward conditional on events in the distant past, since the distant past will have been mostly “forgotten”.

## Learning the $Q$ Function with $n$ -Step SARSA

We might try to learn a  $Q$  function based on partial observations of state by using the obvious generalization of 1-Step SARSA learning:

$$Q(o_t, a_t) \leftarrow (1 - \alpha) Q(o_t, a_t) + \alpha (r_{t+1} + \gamma Q(o_{t+1}, a_{t+1}))$$

But we can't expect this to work, in general —  $Q(o_{t+1}, a_{t+1})$  is *not* the expected discounted future reward from taking  $a_{t+1}$  with observation  $o_{t+1}$  *conditional* on having taken action  $a_t$  the previous time step, when the observation was  $o_t$ .

However, if our policy is ergodic, we should get approximately correct results using  $n$ -Step SARSA for sufficiently large  $n$ . This update for  $Q(o_t, a_t)$  uses actual rewards until enough time has passed that  $a_t$  and  $o_t$  have been (mostly) forgotten:

$$Q(o_t, a_t) \leftarrow (1 - \alpha) Q(o_t, a_t) + \alpha (r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n Q(o_{t+n}, a_{t+n}))$$

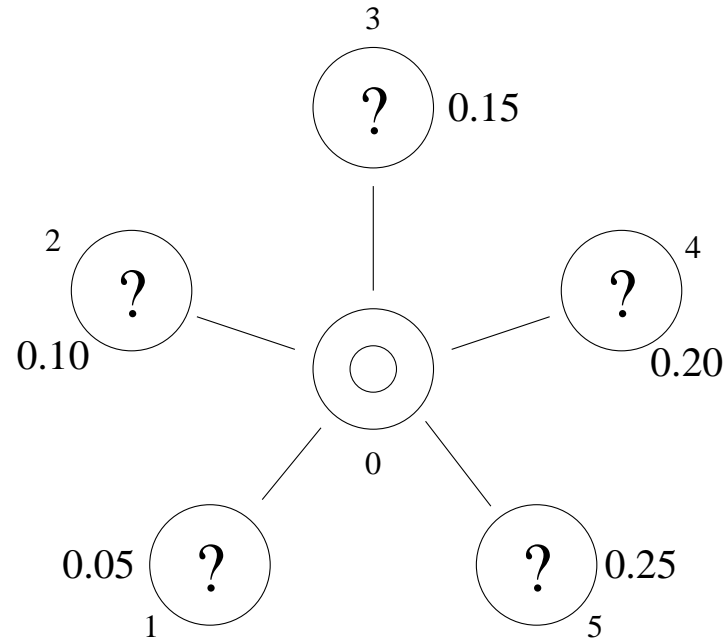
Of course, we have to delay this update  $n$  time steps from when action  $a_t$  was done.

Note!  $n$ -Step SARSA is not the same as SARSA( $\lambda$ ), which can be seen as a weighted combination of  $n$ -Step SARSA for  $n = 1, 2, \dots$  with weights  $1, \lambda, \lambda^2, \dots$

Putting any weight on small  $n$  seems inappropriate here.

# Star World: What Will $Q$ for an “Optimal” Policy Look Like?

Here’s the star world, with the animal in the centre. It can’t see which other locations have food:



Suppose that the animal has no memory of past observations and actions.

What should it do here at the centre, and when at one of the outer locations?

What will the  $Q$  function be like for this policy?



## The Optimal Policy and $Q$ Function

In the star world, we see that without memory, a good policy must be stochastic — sometimes selecting an action randomly.

We can also see that the values of  $Q(o, a)$  for all actions,  $a$ , that are selected with non-zero probability when the observation is  $o$  must be *equal*.

But the probabilities for choosing these actions need not be equal.

So the  $Q$  function for a good policy is not enough to determine this policy.

## But What Does “Optimal” Mean?

But I haven't said what “optimal” means when the state is partially observed. What should we be optimizing?

The most obvious possibility is the average discounted future reward, averaging over the equilibrium distribution of observations (and underlying states):

$$\sum_o P^\pi(o) \sum_a P^\pi(a|o) Q(o, a)$$

Note that the equilibrium distribution of observations depends on the policy being followed, as does the distribution of state given observation.

But with this objective, the discount rate,  $\gamma$ , turns out not to matter!

But it seems to be the most commonly used objective, equivalent to optimizing the long-run average reward per time step.

## But What If We Like Discounted Rewards?

Discounting seems like it's fundamental to decision-making, so this is unsatisfying.

The problem is a conflict between optimizing expected discounted future reward starting from a time when  $o$  is observed versus when  $o'$  is observed:

- We'd like to change  $\pi$  to increase expected discounted reward starting at  $o$ .
- But this could change  $P^\pi(s|o')$  in a way that is bad when  $o'$  is observed later.
- Due to discounting, the bad effect on reward at a later time when  $o'$  is observed is not given full weight when  $o$  is observed.

Proposal: Treat this as a non-cooperative game — the “players” being different observations,  $o$ , and the “moves” being  $P(a|o)$ .

Question: Is there a Nash equilibrium for this game that doesn't require mixed strategies?

I'll assume there is, so we can “optimize” the policy with different criteria for different observations, and still reach an equilibrium.

## Learning a $Q$ Function and an $A$ Function

Since  $Q$  for an optimal stochastic policy does not determine the policy, we can try learning the policy separately, with a similar  $A$  function, updated based on  $Q$ , which is learned with  $n$ -Step SARSA.

The algorithm does the following at each time  $t + n$ :

$$Q(o_t, a_t) \leftarrow (1 - \alpha) Q(o_t, a_t) + \alpha (r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n Q(o_{t+n}, a_{t+n}))$$
$$A \leftarrow A + fQ$$

The policy followed is determined by  $A$ :

$$P(a_t = a \mid o_t = o) \propto \exp(A(o, a) / T)$$

Above,  $T$  is a positive “temperature” parameter, and  $\alpha$  and  $f$  are tuning parameters slightly greater than zero.

This is in the class of “Actor-Critic” methods.

## Learning $A$ , Using $Q$ Not So Much

We can also learn  $A$  more directly, using the same estimates of expected discounted future rewards used to update  $Q$ . We need to weight the updates to  $A$  inversely by the probability of selecting the action taken.

This algorithm does the following at each time  $t + n$ :

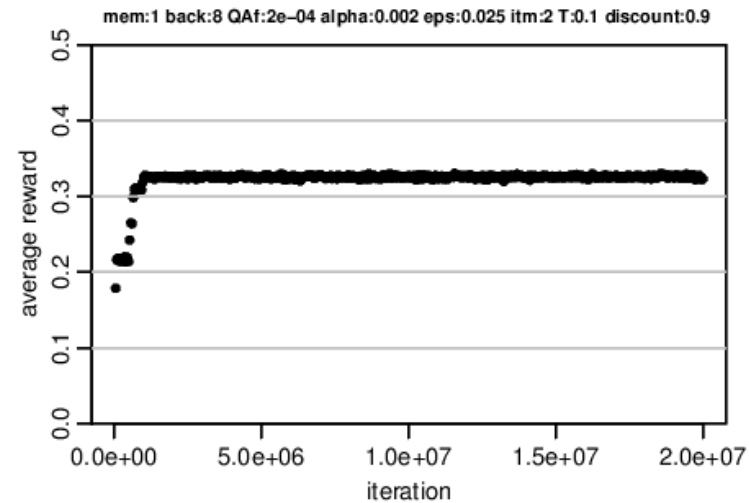
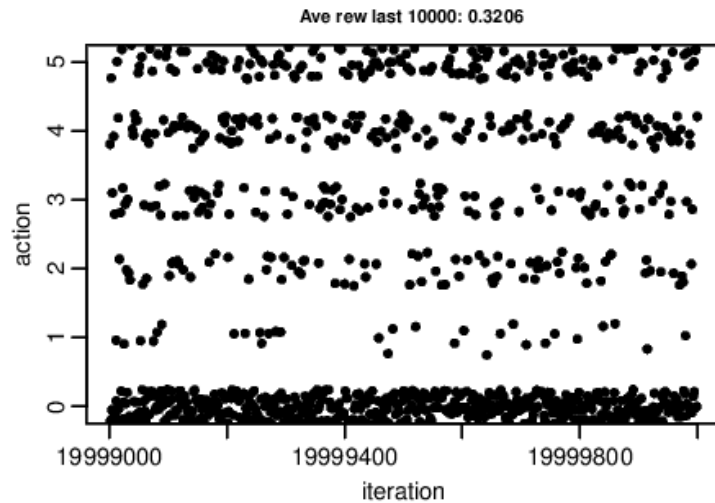
$$E_t \leftarrow r_{t+1} + \gamma r_{t+2} + \cdots + \gamma^{n-1} r_{t+n} + \gamma^n Q(o_{t+n}, a_{t+n})$$

$$Q(o_t, a_t) \leftarrow (1 - \alpha) Q(o_t, a_t) + \alpha E_t$$

$$A(o_t, a_t) \leftarrow A(o_t, a_t) + f E_t / P^{\pi_t}(a_t | o_t)$$

We learn  $Q$  as before, but don't directly use it to update  $A$ . But  $Q$  is still used indirectly, in computing  $E_t$ .

# Star World: Learning $Q$ and $A$



Q:

	0	1	2	3	4	5
0	3.164	3.407	3.360	3.404	3.418	3.380
1	3.135	2.928	2.134	2.154	2.146	2.141
2	3.074	2.103	2.937	2.090	2.118	2.159
3	3.069	2.085	2.093	2.977	2.108	2.120
4	3.059	2.056	2.060	2.092	2.962	2.071
5	3.015	2.059	2.079	2.044	2.072	3.026

P action:

	0	1	2	3	4	5
0	0	5	17	19	28	30
1	98	0	0	0	0	0
2	98	0	0	0	0	0
3	98	0	0	0	0	0
4	98	0	0	0	0	0
5	98	0	0	0	0	0

## So are These Methods Better Than $n$ -Step SARSA?

These methods can learn to pick actions randomly from a distribution that is non-uniform, even when the  $Q$  values for these actions are all the same.

Contrast this with simple  $n$ -Step SARSA, where the  $Q$  function is used to pick actions according to

$$P(a_t = a | s_t = s) \propto \exp(Q(s, a) / T)$$

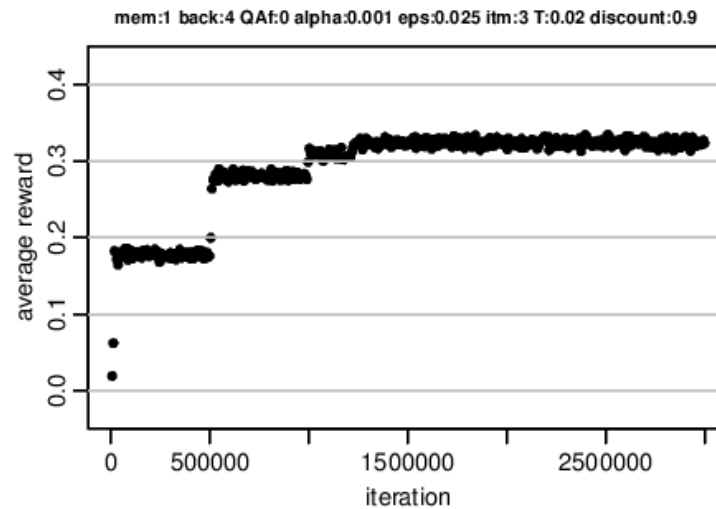
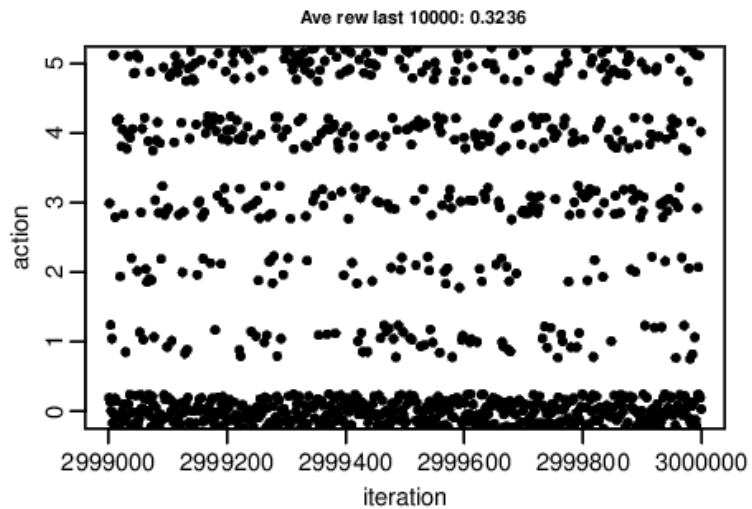
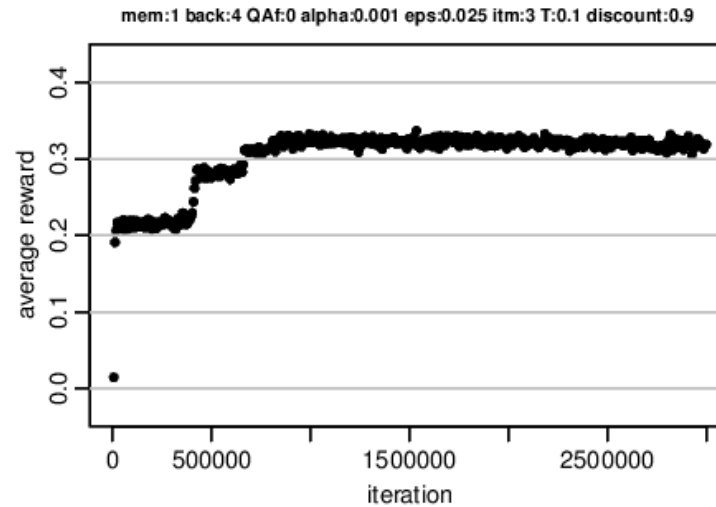
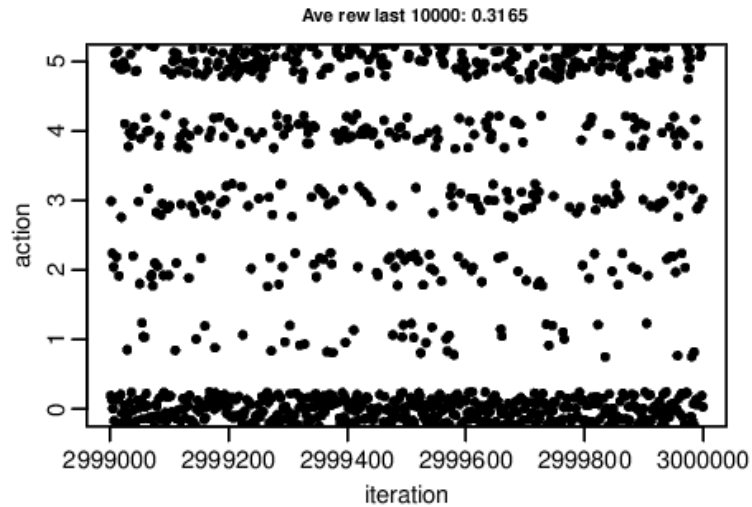
**Obviously**, you can't have  $P(a_t = a | s_t = s) \neq P(a_t = a' | s_t = s)$  when you have  $Q(s, a) = Q(s, a')$ .

Or is it so obvious? What about the limit as  $T$  goes to zero, without being exactly zero?

I figured I should checked it out, just to be sure...

# Using Simple $n$ -Step SARSA With Small $T$ Actually Works!

Here is  $n$ -Step SARSA with  $T = 0.1$  versus  $T = 0.02$ :





# The Policies Learned

The numerical performance difference seems small, but we can also see a qualitative difference in the policies learned:

n-Step SARSA,  $T=0.1$ :

P action:

	0	1	2	3	4	5
0	2	10	15	22	27	24
1	98	0	0	0	0	0
2	98	0	0	0	0	0
3	98	0	0	0	0	0
4	98	0	0	0	0	0
5	60	0	0	0	0	39

n-Step SARSA,  $T=0.02$ :

P action:

	0	1	2	3	4	5
0	0	10	14	25	23	27
1	98	0	0	0	0	0
2	98	0	0	0	0	0
3	98	0	0	0	0	0
4	98	0	0	0	0	0
5	98	0	0	0	0	0

The table entries are probabilities in percent, rounded.

# Comparison of Methods

These methods have different potential deficiencies:

- When learning  $A$  using  $Q$ , we need to learn  $Q$  faster than  $A$ , to avoid changing  $A$  based on the wrong  $Q$ . So  $f$  may have to be rather small (much smaller than  $\alpha$ ).
- When learning  $A$  with inverse probability weights, we may occasionally get a very large weight. This is limited by the probability of exploration,  $\epsilon$ , but still may require a small  $f$ .
- When learning only  $Q$ , with  $T$  very small, the noise in estimating  $Q$  gets amplified by dividing by  $T$ . We may need to make  $\alpha$  small to get less noisy estimates.

- I. Background on Reinforcement Learning with Fully Observed State
- II. Learning Stochastic Policies When the State is Partially Observed
- III. Learning What to Remember of Past Observations and Actions
- IV. Can This Work For More Complex Problems?

# Why and How to Remember

When we can't see the whole state, remembering past observations and actions may be helpful if it helps the agent infer the state.

Such memories could take several forms:

- Fixed memory for the last  $K$  past observations and actions. But  $K$  may have to be quite large, and we'd need to learn how to extract relevant information from this memory.
- Some clever function of past observations — eg, Predictive State Representations.
- Memory in which the agent explicitly decides to record information as part of its actions.

The last has been investigated before (eg, Peshkin, Meuleau, Kaelbling, 1999), but seems to me like it should be investigated more.

# Memories as Observations, Remembering as Acting

We can treat the memory as part of the state, which the agent always observes.

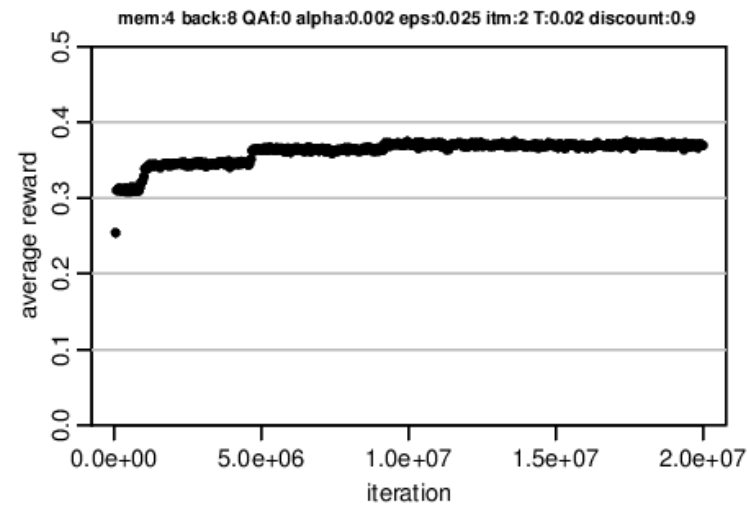
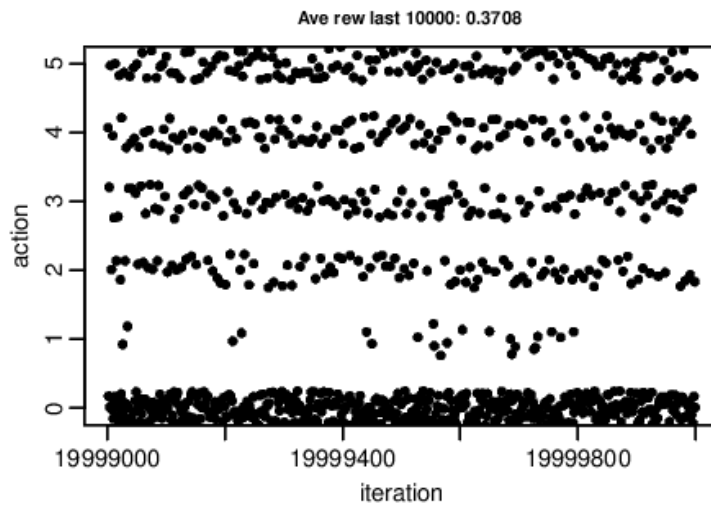
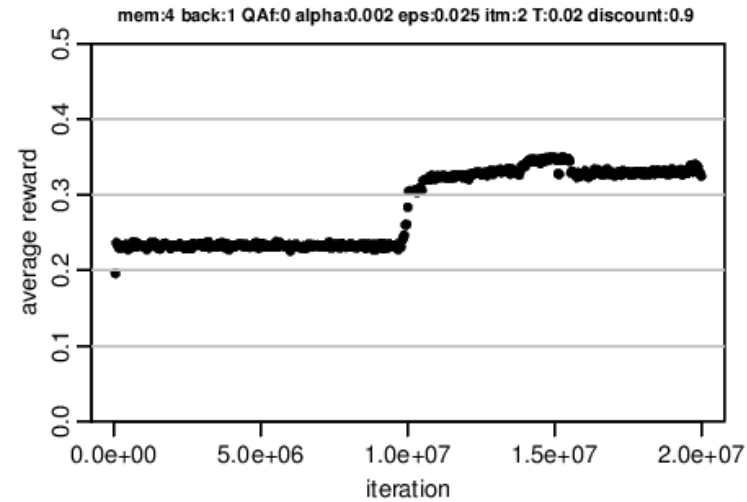
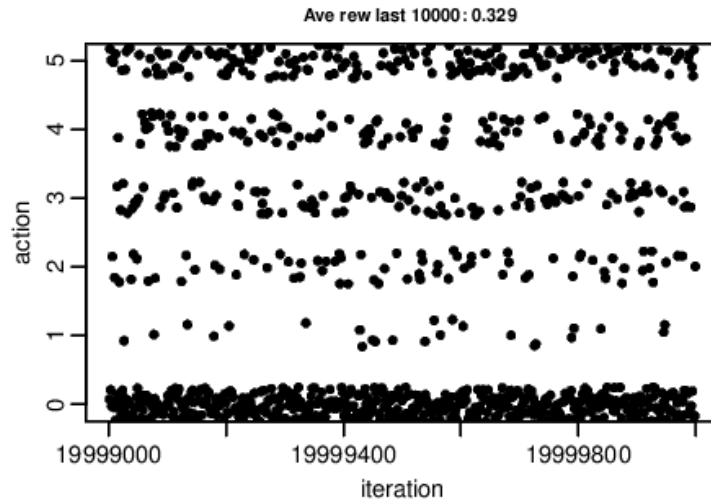
Changes to memory can be treated as part of the action. Most generally, any action could be combined with any change to the memory. But one could consider limiting memory changes (eg, to just a few bits).

Exploration is needed for setting memory as well as for external actions.

In my experiments, I have split exploration into independent exploration of external actions and of internal memory (though both might happen at the same time, with low probability).

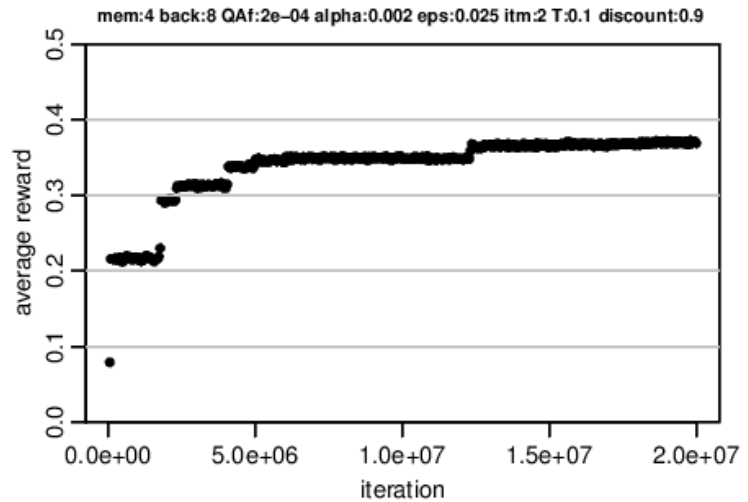
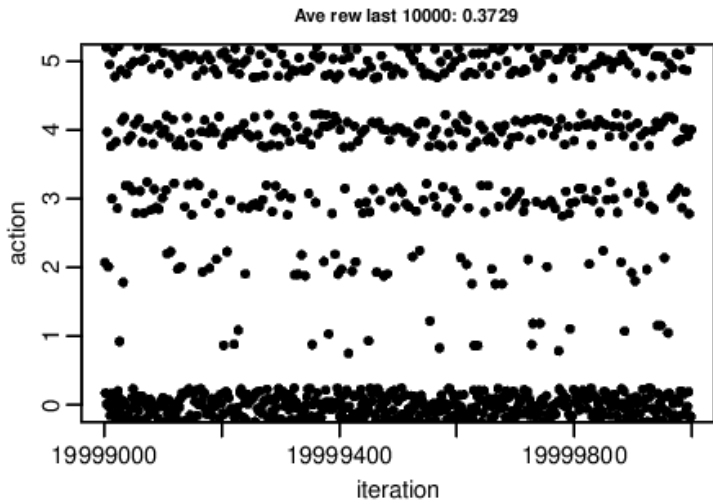
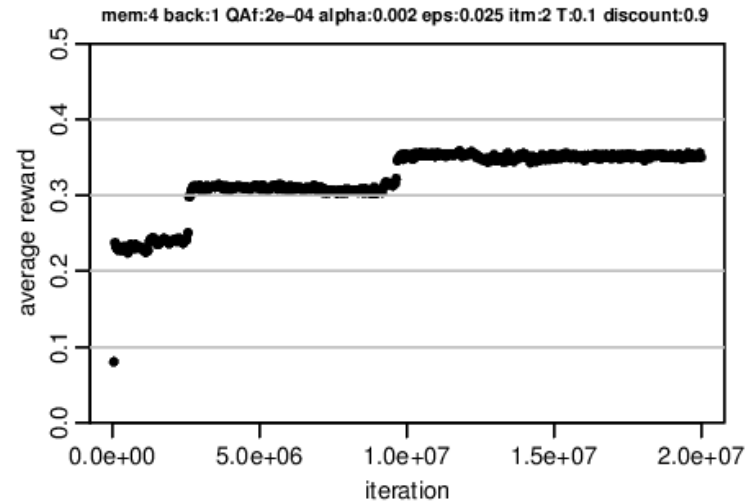
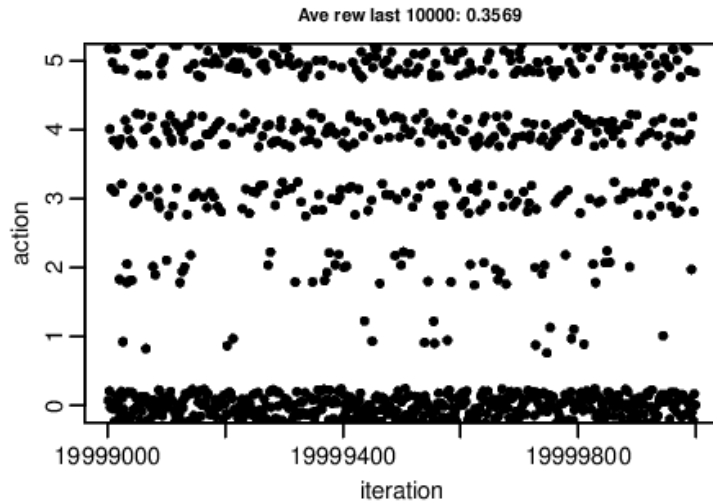
# Star World: 1-Step vs. 8-Step SARSA

4-State Memory, Learns  $Q$



# Star World: 1-Step vs. 8-Step SARSA

4-State Memory, Learns  $Q/A$



- I. Background on Reinforcement Learning with Fully Observed State
- II. Learning Stochastic Policies When the State is Partially Observed
- III. Learning What to Remember of Past Observations and Actions
- IV. Can This Work For More Complex Problems?



# What About More Complex Problems?

As has long been recognized, simple table-based implementations won't work well for complex problems.

Some possibilities I'd like to try:

- Representing  $Q$  and  $A$  Functions with Neural Networks.
- Handling real-valued observations.
- Handling real-valued memory.
- Using ensembles of policies, learned in parallel from the same experiences.

From an AI perspective, I think it's interesting to see how much an agent can learn without detailed guidance — eg, maps of the environment and where the agent is (or may be) located?

## References

Peshkin, L., Meuleau, N., and Kaelbling, L. P. (1999) “Learning Policies with External Memory”, ICML 16.

Singh, S. P., Jaakola, T., and Jordan, M. I. (1994) “Learning without state-estimation in partially observable Markovian decision processes”, ICML 11.