

Reinforcement Learning: Basic models and algorithms

Optimal decisions, Part VII

Christos Dimitrakakis

Chalmers

November 20, 2013

The reinforcement learning problem and MDPs

Markov decision processes

- Description of environments
- Solutions to bandit problems

Algorithms for unknown MDPs.

- Stochastic exact algorithms
- Stochastic estimation algorithms
- Online algorithms

The stochastic n -armed bandit problem

$$\mathcal{P} = \{P_i \mid i = 1, \dots, n\}.$$

$$r_t \mid a_t = i \sim P_i.$$

$$\mathbb{E}_\pi U_t = \mathbb{E}_\pi \sum_{k=t}^T r_k, \quad a_t^* \triangleq \max \{\mathbb{E}(r_t \mid a_t = i) \mid i = 1, \dots, n\}.$$

$$\mathcal{P} = \{P_i(\cdot \mid \omega) \mid \omega \in \Omega\}, \quad r_t \mid a_t = i, \omega^* = \omega \sim P_i(r \mid \omega^*). \quad (1.1)$$

$$\mathbb{E}_\xi^\pi U_t = \mathbb{E}_\xi^\pi \sum_{k=t}^T r_k. \quad (1.2)$$

Algorithm 1 Robbins-Monro bandit algorithm

```

1: input Step-sizes  $(\alpha_t)_t$ , initial estimates  $(\mu_{i,0})_i$ , policy  $\pi$ .
2: for  $t = 1, \dots, T$  do
3:   Take action  $a_t = i$  with probability  $\pi(i \mid a_1, \dots, a_{t-1}, r_1, \dots, r_{t-1})$ .
4:   Observe reward  $r_t$ .
5:    $\mu_{t,i} = \alpha_{i,t} r_t + (1 - \alpha_{i,t}) \mu_{i,t-1}$  // estimation step
6:    $\mu_{t,i} = \mu_{j,t-1}$  for  $j \neq i$ .
7: end for
8: return  $\mu_T$ 

```

Definition 1

ϵ -greedy action selection (w.p. $1 - \epsilon$, select an apparently best action, otherwise a random action)

$$\hat{\pi}_\epsilon^* \triangleq (1 - \epsilon_t) \hat{\pi}_t^* + \epsilon_t \text{Unif}(\mathcal{A}), \quad (1.3)$$

$$\hat{\pi}_t^*(i) = \mathbb{I} \left\{ i \in \hat{\mathcal{A}}_t^* \right\} / |\hat{\mathcal{A}}_t^*|, \quad \hat{\mathcal{A}}_t^* = \arg \max_{i \in \mathcal{A}} \mu_{t,i} \quad (1.4)$$

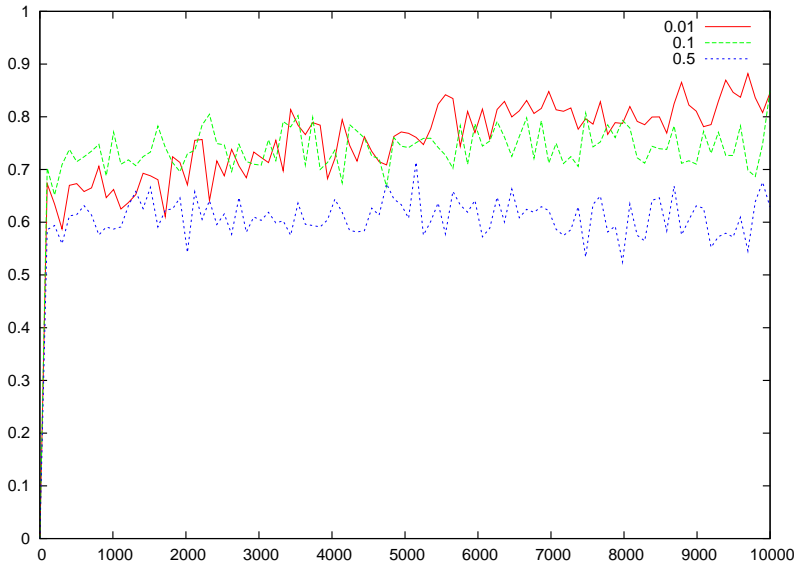


Figure: $\epsilon_t = 0.1$, $\alpha \in \{0.01, 0.1, 0.5\}$.

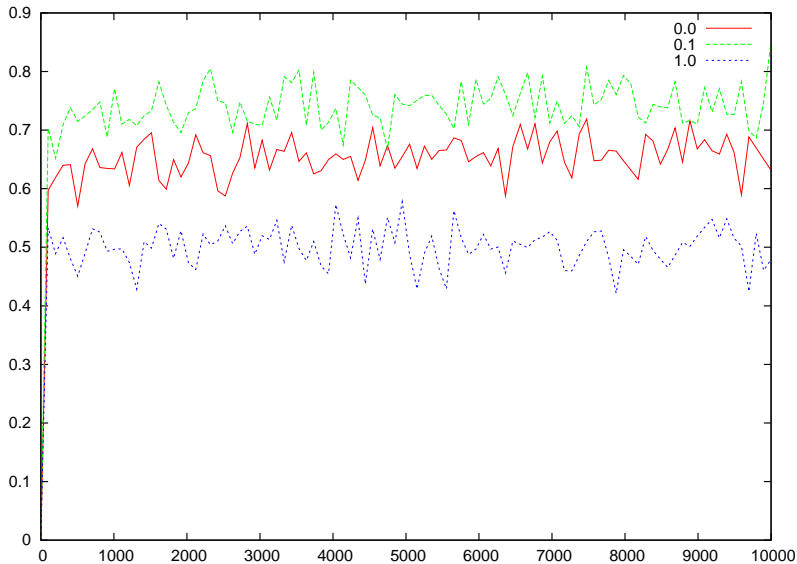


Figure: $\epsilon_t = \epsilon$, $\alpha = 0.1$.

Main idea of the algorithm

- Estimate parameters
- Act according to the estimates

Requirements

- Good estimation procedure.
- Balance estimation with getting rewards!

Consider the algorithm

$$\mu_{t+1} = \mu_t + \alpha_t z_{t+1}. \quad (1.5)$$

Let $h_t = \{\mu_t, z_t, \alpha_t, \dots\}$ be the history of the algorithm.

Assumption 1

Assume a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ such that:

- (i) $f(x) \geq 0$ for all $x \in \mathbb{R}^n$.
- (ii) (Lipschitz derivative) f is continuously differentiable and $\exists L > 0$ such that:

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|, \quad \forall x, y \in \mathbb{R}^n$$

- (iii) (Pseudo-gradient) $\exists c > 0$ such that:

$$c \|\nabla f(\mu_t)\|^2 \leq -\nabla f(\mu_t)^\top \mathbb{E}(z_{t+1} \mid h_t), \quad \forall t.$$

- (iv) $\exists K_1, K_2 > 0$ such that

$$\mathbb{E}(\|z_{t+1}\|^2 \mid h_t) \leq K_1 + K_2 \|\nabla f(\mu_t)\|^2$$

Theorem 2

For the algorithm

$$\mu_{t+1} = \mu_t + \alpha_t Z_{t+1},$$

where $\alpha_t \geq 0$ satisfy

$$\sum_{t=0}^{\infty} \alpha_t = \infty, \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty, \quad (1.6)$$

and under Assumption 1, with probability 1:

- 1 The sequence $\{f(\mu_t)\}$ converges.
- 2 $\lim_{t \rightarrow \infty} \nabla f(\mu_t) = 0$.
- 3 Every limit point μ^* of μ_t satisfies $\nabla f(\mu^*) = 0$.

A demonstration

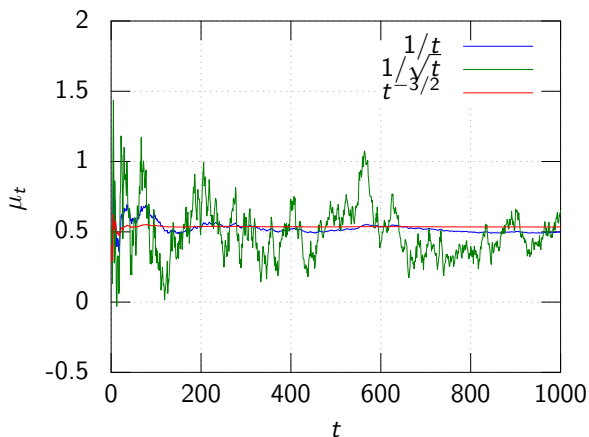


Figure: Estimation of the expectation of $x_t \sim \mathcal{N}(0.5, 1)$ using three step-size schedules.

Algorithm 2 Generic reinforcement learning algorithm

- 1: **input** Update-rule $f : \Theta \times \mathcal{S}^2 \times \mathcal{A} \times \mathcal{R} \rightarrow \Theta$, initial parameters $\theta_0 \in \Theta$, policy $\pi : \mathcal{S} \times \Theta \rightarrow \mathcal{D}(\mathcal{A})$.
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: $a_t \sim \pi(\cdot | \theta_t, s_t)$ // take action
 - 4: Observe reward r_{t+1} , state s_{t+1} .
 - 5: $\theta_{t+1} = f(\theta_t, s_t, a_t, r_{t+1}, s_{t+1})$ // update estimate
 - 6: **end for**
-

Questions

- What should we estimate?
- What policy should we use?

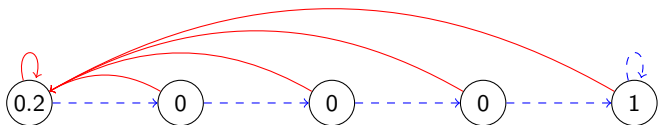


Figure: The chain task

s	s_1	s_2	s_3	s_4	s_5
$V^*(s)$	6.672	7.111	7.689	8.449	9.449
$Q^*(s, 1)$	6.622	6.532	6.676	6.866	7.866
$Q^*(s, 2)$	6.672	7.111	7.689	8.449	9.449

Table: The chain task's value function for $\gamma = 0.95$

Algorithm 3 Stochastic policy evaluation

```

1: input Initial parameters  $v_0$ , Markov policy  $\pi$ .
2: for  $s \in \mathcal{S}$  do
3:    $s_1 = s$ .
4:   for  $k = 1, \dots, K$  do
5:     Run policy  $\pi$  for  $T$  steps.
6:     Observe utility  $U_k = \sum_t r_t$ .
7:     Update estimate  $v_{k+1}(s) = v_k(s) + \alpha_k(U_k - v_k(s))$ 
8:   end for
9: end for
10: return  $v_K$ 

```

For $\alpha_k = 1/k$ and iterating over all \mathcal{S} , this is the same as Monte-Carlo policy evaluation.

Algorithm 4 Approximate policy iteration

```

1: input Initial parameters  $v_0$ , initial Markov policy  $\pi_0$ , stochastic estimator  $f$ .
2: for  $i = 1, \dots, N$  do
3:   Get estimate  $v_i = f(v_{i-1}, \pi_{i-1})$ .
4:   Calculate new policy  $\pi_i = \arg \max_{\pi} \mathcal{L} v_i$ .
5: end for

```

Monte Carlo update

Note that s_1, \dots, s_T contains s_k, \dots, s_T .

Algorithm 5 Every-visit Monte-Carlo update

- 1: **input** Initial parameters v_k , trajectory s_1, \dots, s_T , rewards r_1, \dots, r_T visit counts n .
 - 2: **for** $t = 1, \dots, T$ **do**
 - 3: $U_t = \sum_{s=1}^T r_t$.
 - 4: $n_t(s_t) = n_{t-1}(s_t) + 1$
 - 5: $v_{t+1}(s_t) = v_t(s) + \alpha_{n_t(s_t)}(s_t)(U_t - v_t(s_t))$
 - 6: $n_t(s) = n_{t-1}(s)$, $v_t(s) = v_{t-1}(s) \forall s \neq s_t$.
 - 7: **end for**
 - 8: **return** v_K
-

Example 3

Consider a two-state chain with $\mathbb{P}(s_{t+1} = 1 \mid s_t = 0) = \delta$ and $\mathbb{P}(s_{t+1} = 1 \mid s_t = 1) = 1$, and reward $r(1) = 1$, $r(0) = 0$. Then the every-visit estimate is biased.

Unbiased Monte-Carlo update

Algorithm 6 First-visit Monte-Carlo update

- 1: **input** Initial parameters v_K , trajectory s_1, \dots, s_T , rewards r_1, \dots, r_T , visit counts n .
 - 2: Let $m \in \mathbb{N}^{|S|}$ be trajectory visit counts.
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: $U_t = \sum_{t=1}^T r_t$.
 - 5: $n_t(s_t) = n_{t-1}(s_t) + 1$
 - 6: $m_t(s_t) = m_{t-1}(s_t) + 1$
 - 7: $v_{t+1}(s_t) = v_t(s) + \alpha_{n_t(s_t)}(s_t)(U_t - v_t(s_t))$ if $m_t(s_t) = 1$.
 - 8: $n_t(s) = n_{t-1}(s)$, $v_t(s) = v_{t-1}(s)$ otherwise
 - 9: **end for**
 - 10: **return** v_K
-

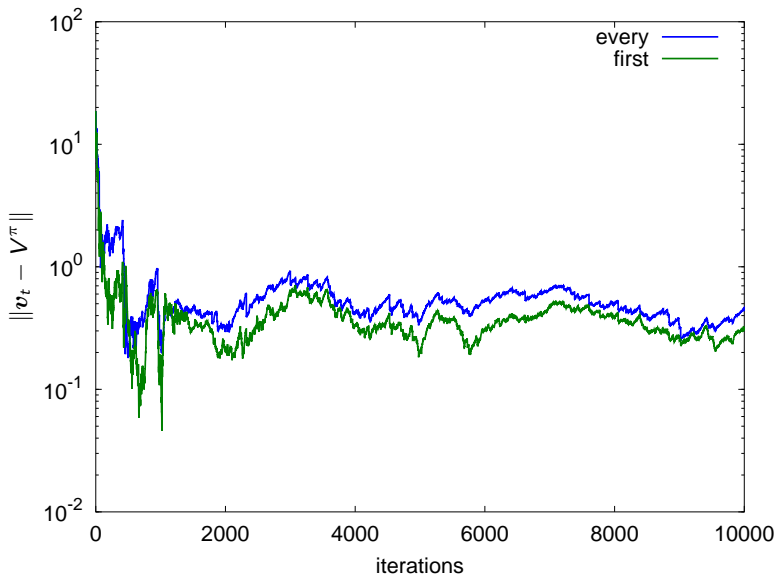


Figure: Error as the number of iterations n increases, for first and every visit Monte Carlo estimation.

Temporal differences

The full stochastic update is of the form:

$$\mathbf{v}_{k+1}(\mathbf{s}) = \mathbf{v}_k(\mathbf{s}) + \alpha(U_k - \mathbf{v}_k(\mathbf{s})),$$

Using the **temporal difference error** $d(\mathbf{s}_t, \mathbf{s}_{t+1}) = \mathbf{v}(\mathbf{s}_t) - [\mathbf{r}(\mathbf{s}_t) + \gamma \mathbf{v}(\mathbf{s}_{t+1})]$,

$$\mathbf{v}_{k+1}(\mathbf{s}) = \mathbf{v}_k(\mathbf{s}) + \alpha \sum_t \gamma^t d_t, \quad d_t \triangleq d(\mathbf{s}_t, \mathbf{s}_{t+1}) \quad (2.1)$$

Stochastic, incremental, update:

$$\mathbf{v}_{t+1}(\mathbf{s}) = \mathbf{v}_t(\mathbf{s}) + \alpha \gamma^t d_t. \quad (2.2)$$

TD(λ)

Temporal-difference operator

$$\mathbf{v}_{n+1}(i) = \mathbf{v}_n(i) + \tau_n(i), \quad \tau_n(i) \triangleq \sum_{t=0}^{\infty} \mathbb{E}_{\pi_n, \mu} [(\gamma \lambda)^m d_n(\mathbf{s}_t, \mathbf{s}_{t+1}) \mid \mathbf{s}_0 = i].$$

Stochastic update:

$$\mathbf{v}_{n+1}(\mathbf{s}_t) = \mathbf{v}_n(\mathbf{s}_t) + \alpha \sum_{k=t}^{\infty} (\gamma \lambda)^{k-t} d_k. \quad (2.3)$$

Algorithm 7 Online TD(λ)

```

1: input Initial parameters  $v_k$ , trajectories  $(s_t, a_t, r_t)$ 
2:  $e_0 = \mathbf{0}$ .
3: for  $t = 1, \dots, T$  do
4:    $d_t \triangleq d(s_t, s_{t+1})$     temporal difference
5:    $e_t(s_t) = e_{t-1}(s_t) + 1$     eligibility increase
6:   for  $s \in \mathcal{S}$  do
7:      $v_{t+1}(s) = v_t(s) + \alpha_t e_t(s) d_t$ .    update all eligible states
8:   end for
9:    $e_{t+1} = \lambda e_t$ 
10: end for
11: return  $v_T$ 

```

Algorithm 8 Simulation-based value iteration

- 1: Input μ, \mathcal{S} .
 - 2: Initialise $s_t \in \mathcal{S}, v_0 \in \mathcal{V}$.
 - 3: **for** $t = 1, 2, \dots$ **do**
 - 4: $s = s_t$.
 - 5: $\pi_t(s) = \arg \max_a \mathbb{P}_\mu(s'|s, a)v_{t-1}(s')$
 - 6: $v_t(s) = r(s) + \sum_{s' \in \mathcal{S}} \mathbb{P}_\mu(s'|s, \pi_t(s))v_{t-1}(s')$
 - 7: $s_{t+1} \sim (1 - \epsilon) \cdot \mathbb{P}(s_{t+1} | s_t = a, \pi_t, \mu) + \epsilon \cdot \text{Unif}(\mathcal{S})$.
 - 8: **end for**
 - 9: Return π_n, V_n .
-

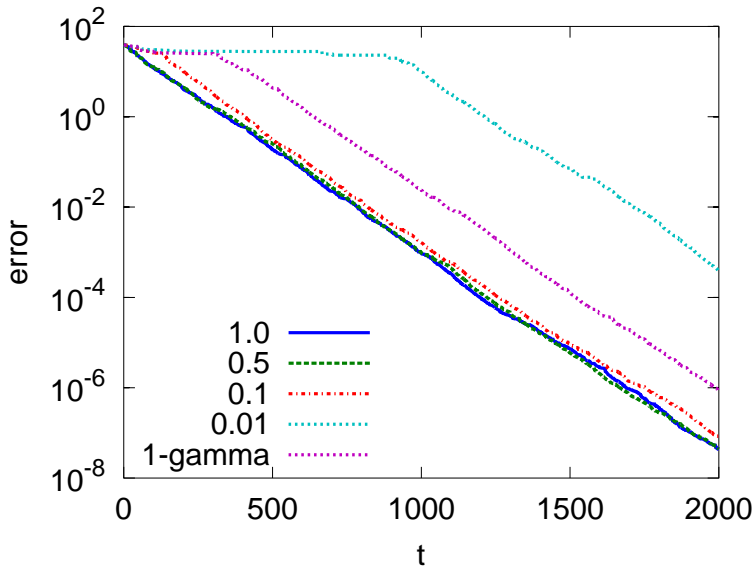


Figure: Simulation-based value iteration with $v_0 = 0$, varying $\epsilon_t = 0.1$.

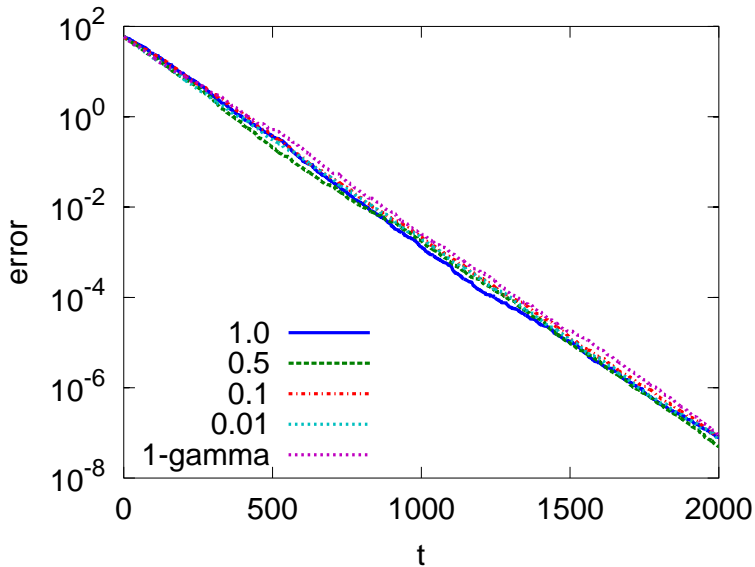


Figure: Simulation-based value iteration with $v_0 = 20 = 1/(1 - \gamma)$, varying ϵ .

Algorithm 9 Q-learning

- 1: Input $\mu, \mathcal{S}, \epsilon_t, \alpha_t$.
 - 2: Initialise $s_t \in \mathcal{S}, \mathbf{q}_0 \in \mathcal{V}$.
 - 3: **for** $t = 1, 2, \dots$ **do**
 - 4: $s = s_t$.
 - 5: $a_t \sim \hat{\pi}_{\epsilon_t}^*(a \mid s_t, \mathbf{q}_t)$
 - 6: $s_{t+1} \sim \mathbb{P}(s_{t+1} \mid s_t = a, \pi_t, \mu)$.
 - 7: $\mathbf{q}_{t+1}(s_t, a_t) = \mathbf{q}_t(s_t, a_t) - \alpha_t[r(s) + v_t(s_{t+1})]$, where $v_t(s) = \max_{a \in \mathcal{A}} \mathbf{q}_t(s, a)$.
 - 8: **end for**
 - 9: Return π_n, V_n .
-

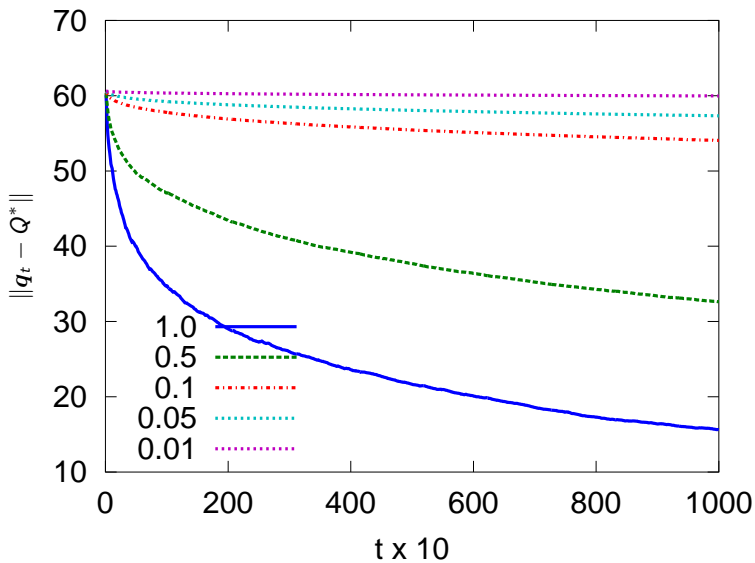


Figure: Q-learning estimation error with $v_0 = 1/(1 - \gamma)$, $\epsilon_t = 1/n_{s_t}$, $\alpha_t \in \alpha n_{s_t}^{-2/3}$.

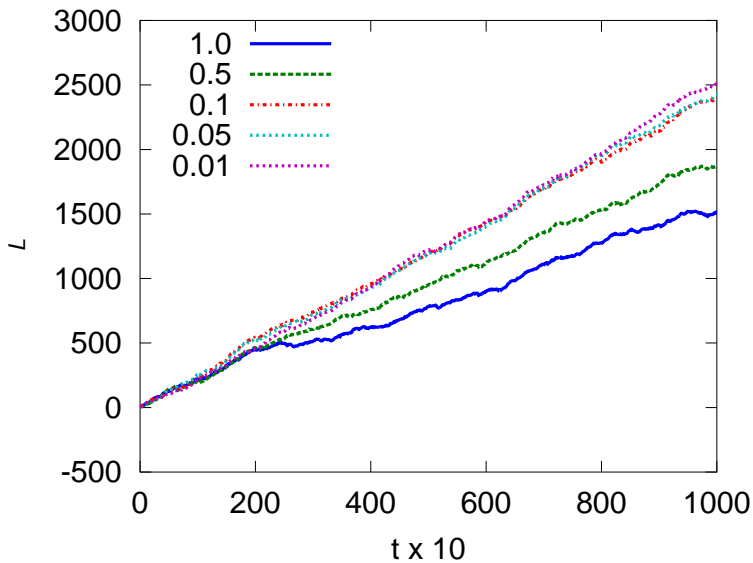


Figure: Q-learning estimation error with $v_0 = 1/(1 - \gamma)$, $\epsilon_t = 1/n_{s_t}$, $\alpha_t \in \alpha n_{s_t}^{-2/3}$.

Algorithm 10 Generalised stochastic value iteration

- 1: Input $\hat{\mu}_0, \mathcal{S}, \epsilon_t, \alpha_t$.
- 2: Initialise $s_1 \in \mathcal{S}, \mathbf{q}_1 \in \mathcal{V}$.
- 3: **for** $t = 1, 2, \dots$ **do**
- 4: $\mathbf{a}_t \sim \hat{\pi}_{\epsilon_t}^*(\mathbf{a} \mid s_t, \mathbf{q}_t)$
- 5: Observe s_{t+1}, r_{t+1} .
- 6: $\hat{\mu}_t = \hat{\mu}_{t-1} \mid s_t, \mathbf{a}_t, s_{t+1}, r_{t+1}$. // update MDP estimate.
- 7: **for** $s \in \mathcal{S}, a \in \mathcal{A}$ **do**
- 8: With probability $\sigma_t(s, a)$ do:

$$\mathbf{q}_{t+1}(s, a) = \mathbf{q}_t(s, a) - \alpha_t \left[r(s) + \gamma \sum_{s' \in \mathcal{S}} \mathbb{P}(s_{t+1} = s' \mid s_t = s, \mathbf{a}_t = a, \hat{\mu}_t) \mathbf{v}_t(s') \right]$$

- 9: otherwise $\mathbf{q}_{t+1}(s, a) = \mathbf{q}_t(s, a)$.
- 10: **end for**
- 11: **end for**
- 12: Return π_n, V_n .

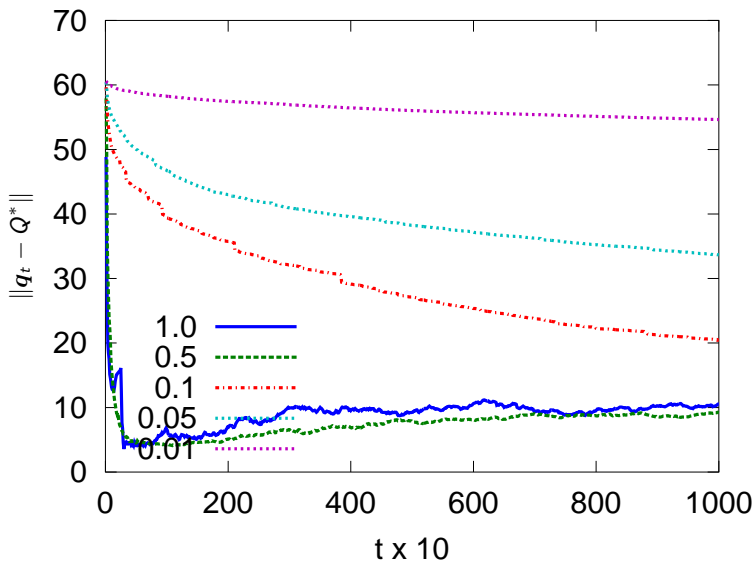


Figure: GSVI with Dirichlet model and single state-action update, with $v_0 = 1/(1 - \gamma)$, $\epsilon_t = 1/n_{s_t}$, $\alpha_t \in \alpha n_{s_t}^{-2/3}$.

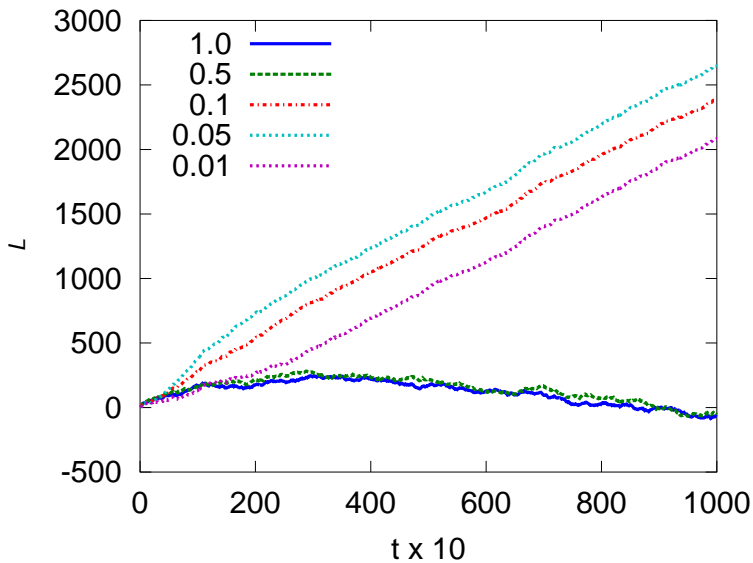


Figure: Q-learning with Dirichlet model and single state-action update with $v_0 = 1/(1 - \gamma)$, $\epsilon_t = 1/n_{s_t}$, $\alpha_t \in \alpha n_{s_t}^{-2/3}$.

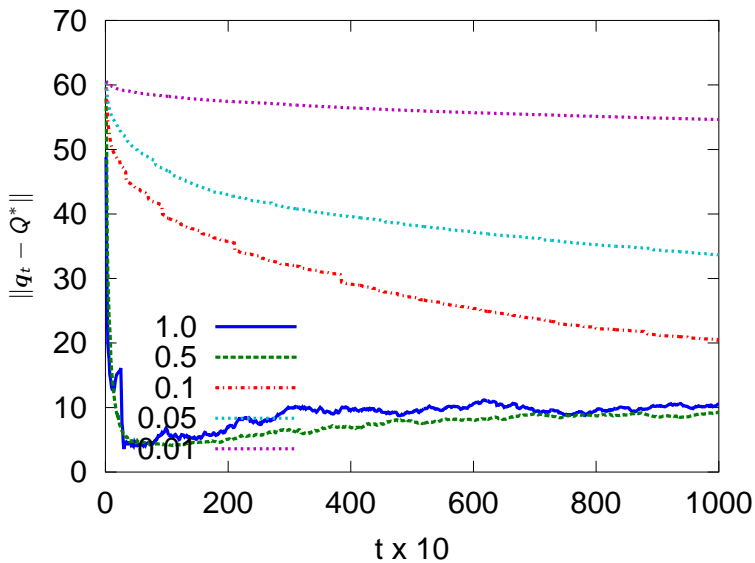


Figure: GSVI with Dirichlet model estimation and a uniform sweep over the state-space with $v_0 = 1/(1 - \gamma)$, $\epsilon_t = 1/n_{s_t}$, $\alpha_t \in \alpha n_{s_t}^{-2/3}$.

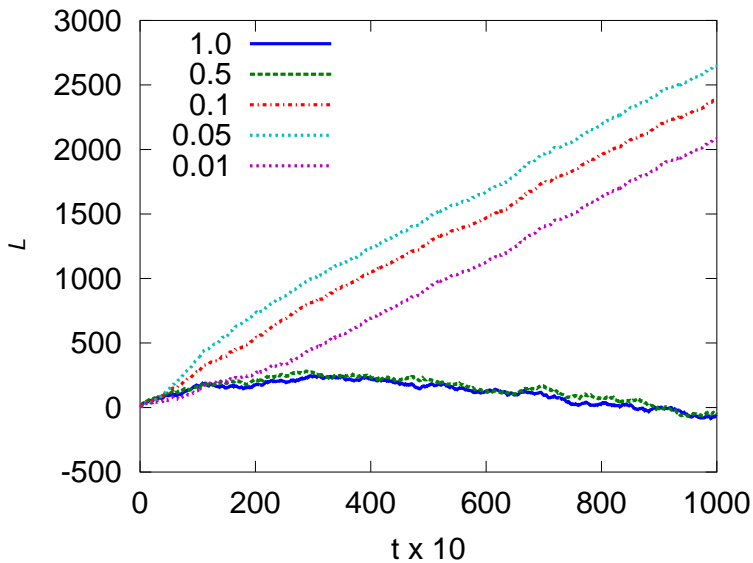


Figure: GSVI with Dirichlet estimation, and a uniform sweep over the state space, with $v_0 = 1/(1 - \gamma)$, $\epsilon_t = 1/n_{s_t}$, $\alpha_t \in \alpha n_{s_t}^{-2/3}$.