

On-policy Control with Approximation

Episodic Semi-gradient Control

- Parameterized action-value function with update target approximation being the Monte Carlo return or any n-step Sarsa returns, for one-step Sarsa

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha \left[R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}_t) - \hat{q}(S_t, A_t, \mathbf{w}_t) \right] \nabla \hat{q}(S_t, A_t, \mathbf{w}_t).$$

- Policy improvement: ε -greedy

Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step size $\alpha > 0$, small $\varepsilon > 0$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Loop for each episode:

$S, A \leftarrow$ initial state and action of episode (e.g., ε -greedy)

 Loop for each step of episode:

 Take action A , observe R, S'

 If S' is terminal:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

 Go to next episode

 Choose A' as a function of $\hat{q}(S', \cdot, \mathbf{w})$ (e.g., ε -greedy)

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

$S \leftarrow S'$

$A \leftarrow A'$

Semi-gradient n-step Sarsa

- N-step return as the update target

$$G_{t:t+n} \doteq R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{q}(S_{t+n}, A_{t+n}, \mathbf{w}_{t+n-1}), \quad t+n < T, \quad (10.4)$$

with $G_{t:t+n} \doteq G_t$ if $t+n \geq T$, as usual. The n -step update equation is

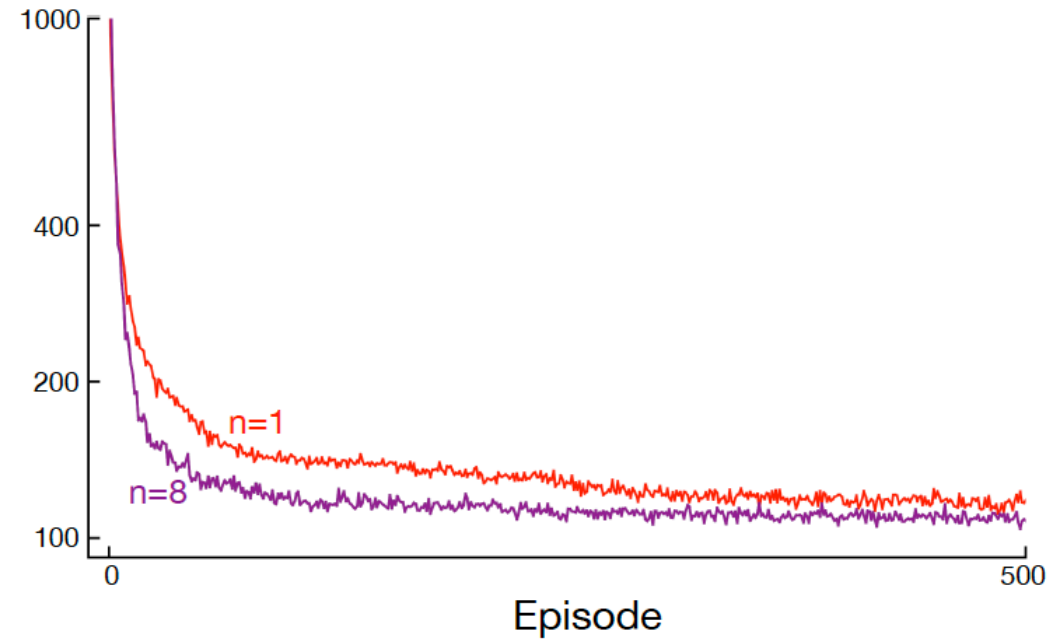
$$\mathbf{w}_{t+n} \doteq \mathbf{w}_{t+n-1} + \alpha [G_{t:t+n} - \hat{q}(S_t, A_t, \mathbf{w}_{t+n-1})] \nabla \hat{q}(S_t, A_t, \mathbf{w}_{t+n-1}), \quad 0 \leq t < T. \quad (10.5)$$

Mountain car

- Grid tilings used to convert the two continuous state variables (position and velocity) to binary features (linear combination)
- Exploration from optimistic action value initialization

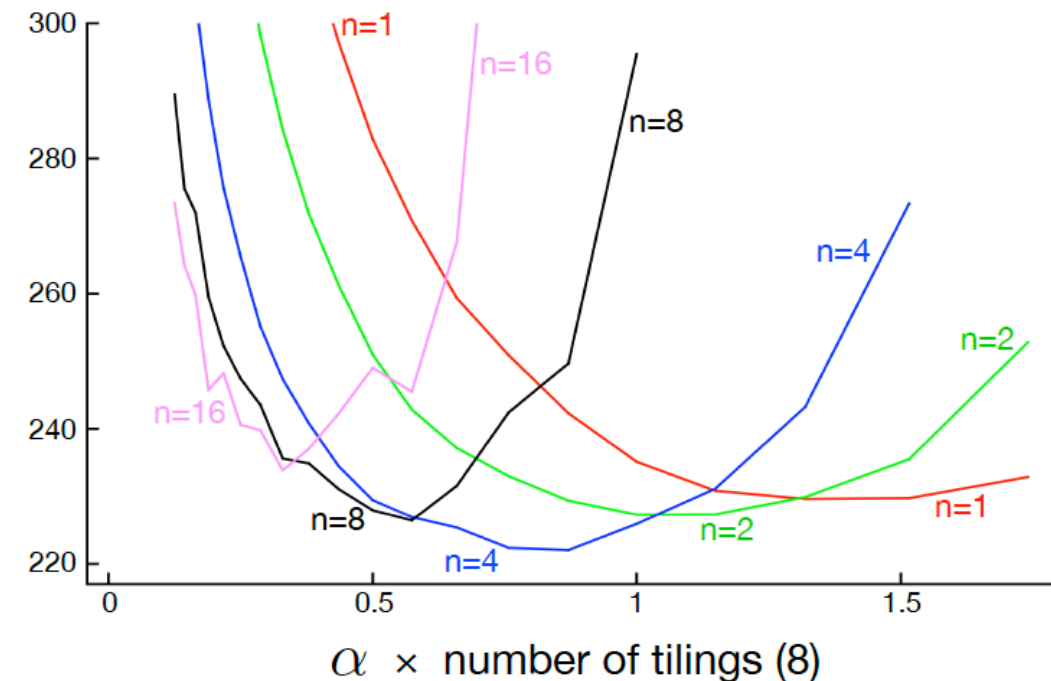
Mountain Car

Steps per episode
log scale
averaged over 100 runs



Mountain Car

Steps per episode
averaged over
first 50 episodes
and 100 runs



Average Reward: A New Problem Setting for Continuing Tasks

- No discounting (commonly considered in the classical theory of DP)

- Average reward following a policy
 - Steady state distribution μ_π assumed to exist for any π and is independent of S_0
= ergodicity

$$\begin{aligned} r(\pi) &\doteq \lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h \mathbb{E}[R_t \mid S_0, A_{0:t-1} \sim \pi] \\ &= \lim_{t \rightarrow \infty} \mathbb{E}[R_t \mid S_0, A_{0:t-1} \sim \pi], \\ &= \sum_s \mu_\pi(s) \sum_a \pi(a|s) \sum_{s', r} p(s', r \mid s, a) r, \end{aligned}$$

- Differential returns: differences between rewards and the average reward

$$G_t \doteq R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \dots$$

Differential value functions

- Differential value functions and Bellman equations
$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{r,s'} p(s', r | s, a) \left[r - r(\pi) + v_{\pi}(s') \right],$$
$$q_{\pi}(s, a) = \sum_{r,s'} p(s', r | s, a) \left[r - r(\pi) + \sum_{a'} \pi(a'|s') q_{\pi}(s', a') \right],$$
$$v_*(s) = \max_a \sum_{r,s'} p(s', r | s, a) \left[r - \max_{\pi} r(\pi) + v_*(s') \right], \text{ and}$$
$$q_*(s, a) = \sum_{r,s'} p(s', r | s, a) \left[r - \max_{\pi} r(\pi) + \max_{a'} q_*(s', a') \right]$$
- Differential TD errors
$$\delta_t \doteq R_{t+1} - \bar{R}_t + \hat{v}(S_{t+1}, \mathbf{w}_t) - \hat{v}(S_t, \mathbf{w}_t),$$
and
$$\delta_t \doteq R_{t+1} - \bar{R}_t + \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}_t) - \hat{q}(S_t, A_t, \mathbf{w}_t),$$

Differential semi-gradient Sarsa for estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step sizes $\alpha, \beta > 0$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Initialize average reward estimate $\bar{R} \in \mathbb{R}$ arbitrarily (e.g., $\bar{R} = 0$)

Initialize state S , and action A

Loop for each step:

 Take action A , observe R, S'

 Choose A' as a function of $\hat{q}(S', \cdot, \mathbf{w})$ (e.g., ε -greedy)

$\delta \leftarrow R - \bar{R} + \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})$

$\bar{R} \leftarrow \bar{R} + \beta \delta$

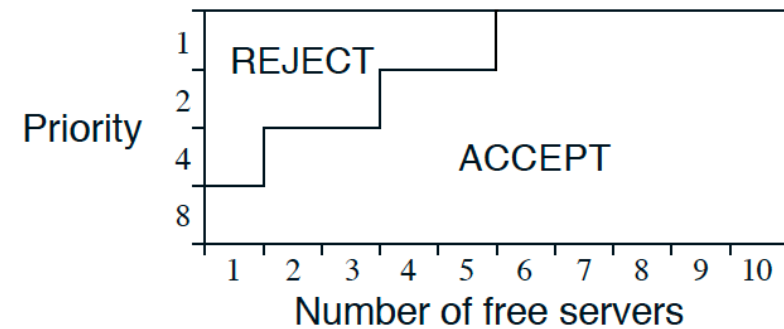
$\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \nabla \hat{q}(S, A, \mathbf{w})$

$S \leftarrow S'$

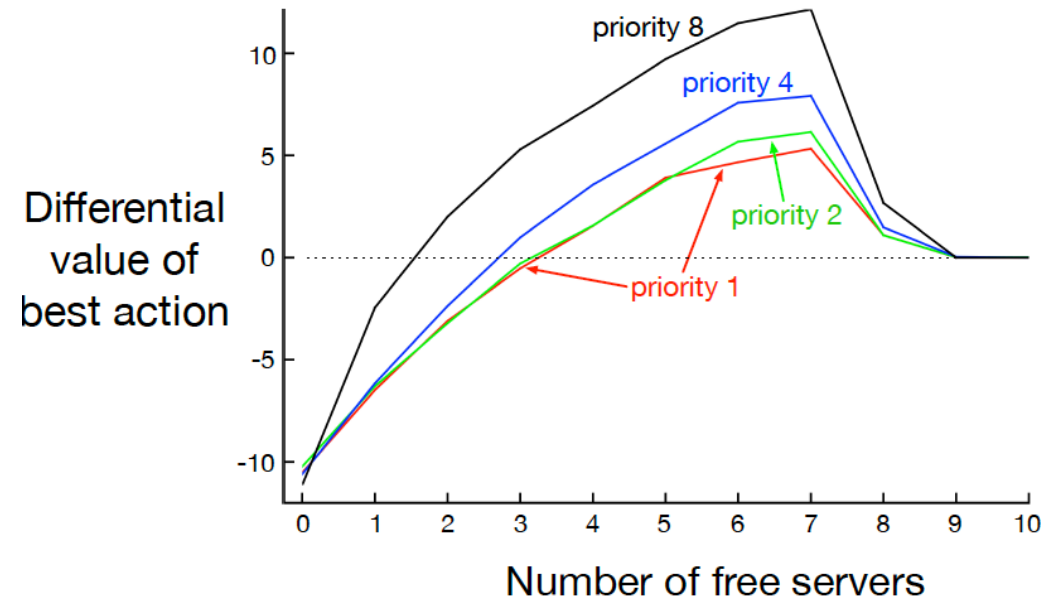
$A \leftarrow A'$

An Access-Control Queuing Task example

- 10 servers, 4 priorities (equal to reward), action = accept or reject head of the queue (never ending with random order), each server becomes free with probability 0.06
- Differential semi-gradient one-step Sarsa



POLICY



VALUE FUNCTION

Deprecating the Discounted Setting

- For continuous control tasks with function approximation, the average of the discounted returns is: $r(\pi)/(1 - \gamma)$,
 - I.e. proportional to the average reward (with the same ordering of all policies)
- γ has no effect due to symmetry
 - Each time step is exactly the same as every other and every reward will appear exactly once in each position in some return
- For control with FA we have lost the policy improvement theorem