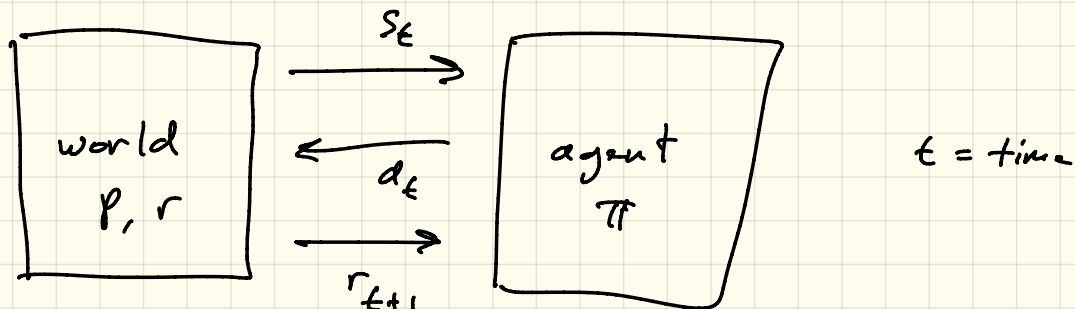# Value - function methods

MDPs — Markov Decision Process



$t$ = time

Spaces: $s_t$ state, $a_t$ action, $r_{t+1}$ reward $\in \mathbb{R}$
$\in \{1, \ldots, N\}'$ $\in \{1, \ldots, M\}'$

$p: S \times S \times A \rightarrow [0, 1]$

functions: model $p(s_{t+1} \mid s_t, a_t)$ = prob. next state is $s_{t+1}$, given $s_t, a_t$

$\pi: A \times S \rightarrow [0, 1]$

policy $\pi(a_t \mid s_t)$ = prob. of taking action $a_t$ in state $s_t$

Sutton & Barto

reward $R(s_t, a_t)$ = reward for action $a_t$ in state $s_t$

<u>notation</u>

$$p(s' \mid s, a) = P(S_{t+1} = s' \mid S_t = s, A_t = a)$$

R.V.s.    $S_{t+1}, S_t, A_t$

<u>Ex</u> :    $f_1(s)$ = expected (avg) 1-step reward starting from $s$

$$= \sum_a \pi(a \mid s) R(s, a)$$

$$= E_{a \sim \pi(\cdot \mid s)} [R(s, a)]$$

$\hookleftarrow$ sampled from

$$= E[R(s, a) \mid s]$$

$$= E[R(S_t, A_t) \mid S_t = s]$$

## Value functions

discount factor    policy    initial state

$$V_\pi(s) = E\left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid \pi, s_0 = s\right]$$

$$Q_\pi(s,a) = E\left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} \mid \pi, s_0 = s, a_0 = a\right]$$

Relationships:

$$V_\pi(s) = E_{a \sim \pi(\cdot|s)}\left[Q_\pi(s,a)\right]$$

$$Q_\pi(s,a) = R(s,a) + \gamma E_{s' \sim p(\cdot|s,a)}\left[V_\pi(s')\right]$$

$$s_t \xrightarrow[a_t]{r_{t+1}} s_{t+1}$$

$$Q(s_t, a_t) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots$$

$$R(s,a)$$

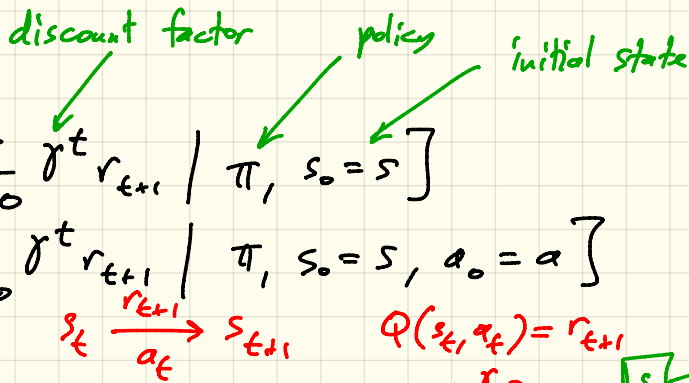$$Q_\pi(s,a) = \underbrace{r_1}_{R(s,a)} + \underbrace{r_2 + r_3 + r_4 + \cdots}_{V_\pi(s')}$$

Optimality:

$$V^*(s) = \max_\pi V_\pi(s)$$

$$Q^*(s,a) = \max_\pi Q_\pi(s,a)$$

agent $\longrightarrow$ best action maximizes $Q^*(s,a)$

$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \arg\max_{a'} Q^*(s,a') \\ 0 & \text{otherwise} \end{cases}$$

RL — learning from sampled action/reward/states

## SARSA

Given $s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}$

$$Q(s_t, a_t) = r_{t+1} + \gamma E_{s_{t+1} \sim p(\cdot|s_t, a_t),\ a_{t+1} \sim \pi(\cdot|s_{t+1})} \left[ Q(s_{t+1}, a_{t+1}) \right]$$

assume $s_{t+1}, a_{t+1}$ are sampled as above

$$Q^+(s_t, a_t) \approx r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) \quad \longleftarrow \text{ correct on average}$$

$$Q(s_t, a_t) \leftarrow (1-\alpha) Q(s_t, a_t) + \alpha Q^+(s_t, a_t) \qquad \alpha = \text{learning rate}$$

$$\leftarrow Q(s_t, a_t) + \alpha \left( \underbrace{Q^+(s_t, a_t) - Q(s_t, a_t)}_{\delta_t\ \text{target}} \right) \qquad \in [0, 1]$$

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

how do we sample $a_{t+1}$ ?

$$a_{t+1} = \arg\max_a Q(s_{t+1}, a)$$

← greedy policy

pure exploitation

this means that SARSA is "on-policy"

a is determined by current Q

practize: use <u>ε-greedy policy</u> :

$$a_{t+1} = \begin{cases} \text{random } a & \text{with prob. } \varepsilon \\ \arg\max_a Q(s_{t+1}, a) & \text{otherwise} \end{cases}$$
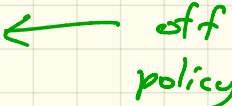
$\varepsilon$ = prob of exploring $\in [0, 1]$

thm : SARSA converges if all $(s, a)$ pairs are visited infinitely often and policy converges to greedy.

For example, anneal $\varepsilon = 1/t$

# Q-learning

Given $s_t, a_t, r_{t+1}, s_{t+1}$

$$Q(s_t, a_t) = r_{t+1} + \gamma E_{s_{t+1} \sim p(\cdot | s_t, a_t), \ a_{t+1} \sim \pi(\cdot | s_{t+1})} \left[ Q(s_{t+1}, a_{t+1}) \right]$$

Do not assume that $a_t / a_{t+1}$ are from $Q$ $\longleftarrow$ off policy

Then $\quad Q(s_t, a_t) = r_{t+1} + \gamma E_{s_{t+1}} \left[ \max_{a'} Q(s_{t+1}, a') \right]$

$$Q^\dagger(s_t, a_t) = r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') \qquad s_{t+1} \text{ is sampled}$$

Update $\quad Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \underbrace{\left( r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)}_{\delta_t}$

In practice, choose $a_{t+1}$ from $\varepsilon$-greedy policy for $Q$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)$$

$$\underbrace{\phantom{r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)}}_{\text{target}} \qquad \nabla L ?$$

$$L_t(q) = \frac{1}{2} \left( \underbrace{r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a')}_{\text{target}} - q \right)^2$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) - \alpha \nabla_q L_t \left( Q(s_t, a_t) \right)$$