

Lecture 2 : Tabular RL

Last time: Value and Policy Iteration

Value Iteration

$$Q_0 = 0.$$

For  $t = 0, 1, \dots$

$$Q_{t+1} = TQ_t$$

————— x ————— x ————— x —————

Policy Iteration

Initial policy  $\mu_0$

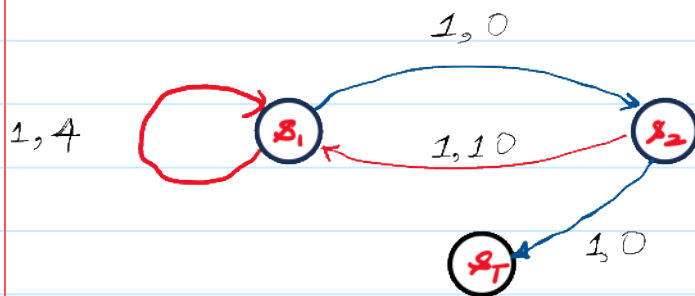
For  $t = 0, 1, \dots$

Policy Evaluation

Find  $Q_{\mu_t}$

Policy Update

Set  $\mu_{t+1} = \text{greedy}(Q_{\mu_t})$

Illustrative Example

$$\text{Let } Q_0 = [0, 0, 0, 0, 0]^T.$$

$$\& \text{ } r = 0.9.$$

$$\text{Then, } Q_1 = T Q_0$$

$$\Rightarrow Q_1(s_1, L)$$

$$= u(s_1, L) + r \max_{a' \in \{L, R\}} Q_0(s_1, a')$$

$$= 4.$$

Similarly,

$$Q_1(s_1, R) = 0$$

$$+ r \max_{a' \in \{L, R\}} Q_0(s_2, a')$$

$$= 0.$$

Continuing this way, we get

$$Q_1(s_0, L) = 10$$

$$Q_1(s_0, R) = 0$$

$$\nexists Q_2(s_T, a) = 0 \quad \forall a \in \{L, R\}.$$

Observe that the greedy policy at  $s_1$   $\nexists$   $s_2$  suggest picking the action L at both places

Let us compute  $Q_2$  now

$$Q_2(s_1, L) = 4 +$$

$$\gamma \max \{Q_1(s_1, L), Q_1(s_1, R)\}$$

$$= 4 + 0.9 \times 4$$

$$= 7.6.$$

$$Q_2(s_1, R) = 0 +$$

$$\gamma \max \{Q_1(s_2, L), Q_1(s_2, R)\}$$

$$= 0.9 \times 10 = 9.$$

$$Q_2(s_2, L)$$

$$= 10 + \gamma \max \left\{ Q_1(s_1, L), Q_1(s_1, R) \right\}$$

$$= 10 + 0.9 \max \{ 4, 0 \}$$

$$= 13.6$$

$$Q_2(s_2, R) = 0$$

————— x ————— x ————— x —————

The behaviour of the Policy Iteration method can similarly be studied for the above example.

————— → ————— x ————— → —————

We now look at model-free variants of the above algorithm.

These algorithms do not need knowledge of the transition probabilities, instead they learn using sampled trajectories.

To begin with, we first look at the problem of policy evaluation.

there, given a policy  $\mu$ , the goal is to find

$Q_\mu$ .

Recall that, for any  $\pi, \mu$

$Q_\mu(\pi, a)$

$$= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R_t(s_t, a_t) \mid s_0 = \pi, a_0 = a \right]$$

and the expectation is with respect to

$$s_{t+1} \sim p(\cdot \mid s_t, a_t)$$

and

$$a_{t+1} \sim \mu(\cdot \mid s_{t+1})$$

for all  $t \geq 0$ .

## Temporal-Difference Learning for Policy Evaluation.

We presume  $\mu(a|s) > 0 \forall s, a$ .

Initialize  $Q_0$  and initial state-action pair  $s_0, a_0$

For  $t \geq 0$ ,

Observe  $r_t$

Observe  $s_{t+1} \sim p(\cdot | s_t, a_t)$

sample  $a_{t+1} \sim \mu(\cdot | s_{t+1})$

Update  $Q_{t+1}(s, a)$

$$= \begin{cases} Q_t(s, a), & \text{if } (s, a) \neq (s_t, a_t) \\ Q_t(s_t, a_t), & \text{otherwise} \end{cases}$$

$$+ \alpha_t [r_t + \gamma Q_t(s_{t+1}, a_{t+1}) - Q_t(s_t, a_t)]$$

$$+ \gamma Q_t(s_{t+1}, a_{t+1})$$

$$- Q_t(s_t, a_t)]$$

where  $(\alpha_t)$  is some step-size sequence.

Remark: This algorithm is known as TD(0) and it updates  $Q_n$  as soon as one transition occurs.

Intuition behind algorithm design

Let

$$f(Q) = \frac{1}{2} \|Q - Q_4\|_{D_4}^2$$

where  $D_4$  is a  $SA \times SA$  diagonal matrix whose  $(s,a)$ -th diagonal entry is

$$d_{s,a} = \mu(a|s)$$

and  $d_{s,a}$  is the stationary distribution of the Markov chain induced by  $\mu$ .

Thus,

$$f(Q) = \frac{1}{2} \sum_{x,a} d_y(x) \psi(a/x) \\ \times (Q(x,a) - Q_y(x,a))^2 \\ = \frac{1}{2} (Q - Q_y)^T D_y (Q - Q_y)$$

Suppose  $\psi(a/x) > 0 \forall x,a$   
 of the Markov chain  
 induced by  $\psi$  is ergodic  
 so that a unique  
 stationary distribution  
 $d_y$  exists &  $d_y(x) > 0 \forall x$ .

Then,  $D_y$  is positive definite &

$$\nabla f(Q) = D_y (Q - Q_y)$$

Hence, if we would  
 have known  $D_y$  &  $Q_y$ ,

then one could have  
 used the update rule

$$Q_{n+1} = Q_n$$

$$+ \alpha_n D_y [Q_y - Q_n]$$

which is the gradient descent method for the optimization problem

$$\min_{\theta} f(\theta).$$

Since  $\theta_y$  is unknown, we now use its properties to modify the update rule.

Recall

$$T_{\gamma} x = x + \gamma P_{\gamma} x,$$

where  $P_{\gamma}$  is  $SA \times SA$  of

$$P_{\gamma}(s', a' | s, a)$$

$$= P(s' | s, a) U(a' | s').$$

Moreover,

$$T_{\gamma} \theta_{\gamma} = \theta_{\gamma}.$$

$$\begin{aligned} \text{Then, } T_{\gamma} \theta_{\gamma} &= x + \gamma P_{\gamma} \theta_{\gamma} \\ &= \theta_{\gamma}. \end{aligned}$$

In particular,

$$D_y Q_y$$

$$= \sum_{s, a} e_{s, a} d_y(x) u(a/x) Q_y(s, a)$$

$$= \sum_{s, a, s', a'} e_{s, a} d_y(x) u(a/x) p(s'/s, a) u(a'/s')$$

$$\left[ u(s, a) + r Q_y(s', a') \right]$$

$$\approx \sum_{s, a, s', a'} e_{s, a} d_y(x) u(a/x)$$

$$p(s'/s, a) u(a'/s')$$

$$\left[ u(s, a) + r Q_n(s', a') \right]$$

$$= E \left[ e_{s, a} \left( u(s, a) + r Q_n(s', a') \right) \mid \mathcal{F}_n \right],$$

where  $\mathcal{F}_n = \sigma(Q_0, s_0, a_0, s_1, a_1, \dots, s_{n-1}, a_{n-1}, s_n, a_n)$

if the expectation is  
presuming

$$z \sim d_y(\cdot)$$

$$a \sim u(\cdot | z)$$

$$z' \sim p(\cdot | z, a)$$

$$a' \sim u(\cdot | z').$$

Hence, the SGD  
can be approximated  
by

$$\begin{aligned} Q_{n+1} = & Q_n + d_n [x(z_n, a_n) \\ & + \gamma Q_n(z_{n+1}, a_{n+1}) \\ & - Q_n(z_n, a_n)] e_{z_n, a_n}. \end{aligned}$$

which is exactly the  
TD(0) method.

Observe that

$$Q_{n+1} = Q_n + \Delta_n D_y [T_y Q_n - Q_n] \\ + \Delta_n M_{n+1},$$

where

$$M_{n+1} = \epsilon_n \epsilon_{n+1} \Delta_n \\ - D_y [T_y Q_n - Q_n]$$

and

$$\epsilon_n = \mu(x_n, a_n) + \sigma Q_n(x_n, a_n) \\ - Q_n(x_n, a_n).$$

Using stochastic approximation theory, it can be shown that the  $(Q_n)$  iterates track the solution trajectories of the ODE

$$\dot{Q}(t) = D_y [T_y Q_n - Q_n]$$

Additionally, by using the fact that

$T_y$  is a contraction,

it can be shown that

$Q_y$  is a globally asymptotically stable equilibrium for this limiting ODE.

Using this, it can be concluded that

$$Q_t \xrightarrow{a.s.} Q_y$$

as  $t \rightarrow \infty$



### Optimal Control:

Using a similar recipe, one can come up with the following algorithm to find  $Q^*$ :

$$Q_{t+1} = Q_t +$$

$$dt [u(x_t, a_t)]$$

$$+ \max_{a_t} Q_t(x_t, a_t)$$

$$- Q_t(x_t, a_t)] e_{\beta_t} dt$$

At time  $t$ ,  $a_t$  can be chosen in multiple ways:

① fixed behaviour policy

$$a_t \sim \mu(\cdot | s_t)$$

where  $\mu$  is some fixed behaviour policy

②  $\epsilon$ -greedy behaviour policy:

$$a_t = \begin{cases} \text{arbitrary action w.p. } \epsilon \\ \arg \max_{a'} Q_t(s_t, a') \end{cases}$$

w.p.  $(1-\epsilon)$

In either case,

$$Q_t \xrightarrow{a^*} Q^*$$

as  $t \rightarrow \infty$ .