

Margarita Prikhodko

$$1. a) q_n(s, a) = E_{\pi} [G_t \mid S_t = s,$$

$$A_t = a]$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$q_n(s, a) = E_{\pi} [G_t \mid S_t = s,$$

$$A_t = a] = E_{\pi} [R_{t+1} + \gamma G_{t+1}]$$

$$S_t = s, A_t = a] = E_{\pi} [R_{t+1} +$$

$$+ \gamma \sum_{s', a'} q_{\pi}(s', a') \mid S_t = s, A_t = a]$$

$$= \sum_{s', r} p(s', r \mid s, a) [r + \gamma \sum_{a'} \pi(a' \mid s) q_{\pi}(s', a)]$$

$$q_{k+1}(s, a) = E_{\pi} [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

$$= \sum_{s', r} p(s', r \mid s, a) [r + \gamma \sum_{a'} \pi(a' \mid s) q_k(s', a)]$$

b)

$$v_{k+1}(s) = E_{\pi} [R_{t+1} + \gamma v_E(s_{t+1}) \mid S_t = s]$$

$$= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_E(s')]$$

$$= T_{\pi}(v_k)$$

$$q_{k+1} = T_{\pi}^{\gamma}(q_k)$$

Show T_π^q is γ -contracting
 where $0 < \gamma < 1$ is discounting factor.

Let's consider 2 state value functions g_1 and g_2

1) g_1 and g_2 2 state-action value function represented as $n \times m$ dim. vectors. We have to show:

$$\|T_\pi^q(g_1) - T_\pi^q(g_2)\| \leq \gamma \|g_1 - g_2\|$$

$\|\cdot\|$ is Euclidean norm i.e. not of x

$$T_\pi^q(g) = R_g^\pi + P_g^\pi \cdot g$$

Difference between g_1, g_2 output:

$$\begin{aligned} T_\pi^q(g_1) - T_\pi^q(g_2) &= \\ &= (R_{g_1}^\pi + P_{g_1}^\pi \cdot g_1) - (R_{g_2}^\pi + P_{g_2}^\pi \cdot g_2) \\ &= P_{g_1}^\pi \cdot g_1 - P_{g_2}^\pi \cdot g_2 \end{aligned}$$

2 output operators!

$$P_{g_1}^\pi, P_{g_2}^\pi$$

$$\| T_{\bar{q}}^{\pi}(g_1) - T_{\bar{q}}^{\pi}(g_2) \| =$$

$$= \| P_{\bar{q}}^{\pi} \cdot g_1 - P_{\bar{q}}^{\pi} \cdot g_2 \|$$

$$\| P_{\bar{q}}^{\pi} \cdot g_1 - P_{\bar{q}}^{\pi} \cdot g_2 \|$$

Matrix normalization

$$\| X \| = \max (\| X \cdot v \| / \| v \|)$$

From above let show
that $\| P_{\bar{q}}^{\pi} \| \leq \gamma$

$$\| P_{\bar{q}}^{\pi} \cdot g_1 - P_{\bar{q}}^{\pi} \cdot g_2 \| \leq \| P_{\bar{q}}^{\pi} (g_1 - g_2) \|$$

$$\| P_{\bar{q}}^{\pi} \cdot (g_1 - g_2) \| \leq \| P_{\bar{q}}^{\pi} \| \cdot \| g_1 - g_2 \|$$

When $\| P_{\bar{q}}^{\pi} \| \leq \gamma$

$$\| P_{\bar{q}}^{\pi} (g_1 - g_2) \| \leq \gamma \| g_1 - g_2 \|$$



$$\| T_{\bar{\pi}}^{\bar{q}}(g_1) - T_{\bar{\pi}}^{\bar{q}}(g_2) \| \leq \gamma \| g_1 - g_2 \|$$

$T_{\bar{\pi}}^{\bar{q}}$ is contracting

C. R_{π}^q - reward for Taking action
action under Policy π

P_q^{π} State transition

probabilities under policy

The property P_q^{π}

of the q operator

satisfies the Bellman

equation for a function

$$Q_{\pi}(s, a) = R_q^{\pi}(s, a) +$$

$$+\gamma \sum P_q^{\pi}(s'|s, a) \cdot \max_{a'} Q_{\pi}(s'a')$$

Equation is Q-value for
state-action pair (s, a)

depends on the immediate
reward and discount
factor γ .

The contraction property
arises from Bellman equation
a-learning iteration

2. Exercise 4.5

1.

Initialize

$Q(s, a) \in \mathbb{R}$ and

$\pi(s) \in A(s)$ arbitrarily
for $s \in S, a \in A$

Evaluation

loop:

$$\Delta \leftarrow 0$$

loop for each $s \in S$ and $a \in A$:

$$q = Q(s, a)$$

$$Q(s, a) \leftarrow \sum_{s', r} p(s', r | s, a) \cdot$$

$$r + \gamma \sum a' \pi(a'|s) Q(s', a')$$

$$\Delta \leftarrow \max(\Delta, |q - Q|)$$

while $\Delta < \epsilon$ (ϵ - is error
small pos. number)

Improvement

policy-stable \leftarrow true

For each $s \in S$ and $a \in A$

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \arg \max_a Q(s, a)$

If old action \neq $\pi(s)$ then

Then policy-stable \leftarrow false

If policy stable
then stop and return
 $Q \approx q^*$ and $\pi \approx \pi^*$
else go to 2.