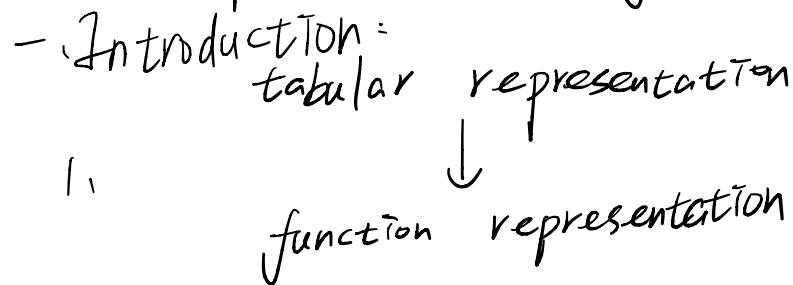



Chapter 8 : value function approximation



So far : State and action value are represented by tables, (表格)

- For example, action value:

$$q_{\pi}(s, a)$$

	a_1	a_2	a_3	a_4	a_5
s_1	$q_{\pi}(s_1, a_1)$	$q_{\pi}(s_1, a_2)$	$q_{\pi}(s_1, a_3)$	$q_{\pi}(s_1, a_4)$	$q_{\pi}(s_1, a_5)$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
s_9	$q_{\pi}(s_9, a_1)$	$q_{\pi}(s_9, a_2)$	$q_{\pi}(s_9, a_3)$	$q_{\pi}(s_9, a_4)$	$q_{\pi}(s_9, a_5)$

① storage

存储能力.

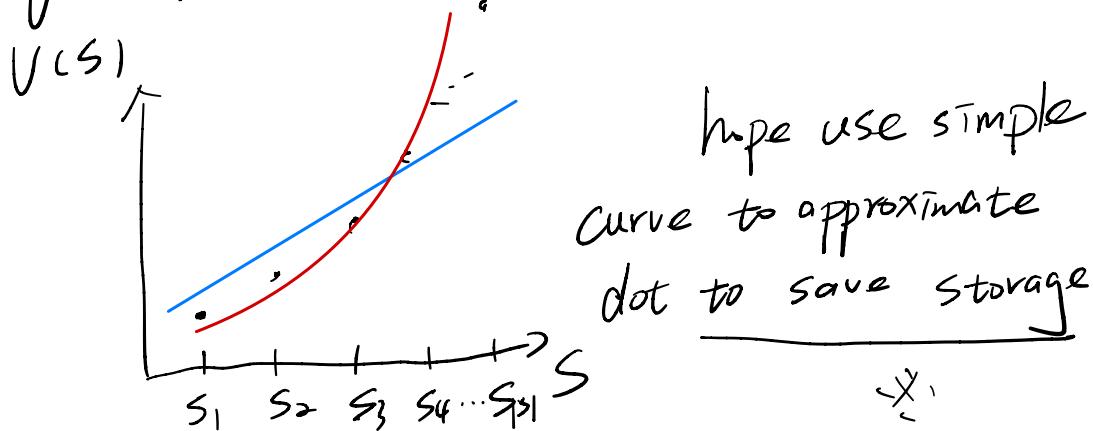
② generalization ability

泛化能力

2. Consider Examples.

One-dimensional state = $s_1, \dots, s_{|S|}$

π is given: State value: $V_\pi(s_1), V_\pi(s_2) \dots V_\pi(s_{|S|})$



$$\text{Suppose straight } V(s, w) = a s + b = \underbrace{[s, 1]}_{\Phi^T(s)} \begin{bmatrix} a \\ b \end{bmatrix} = \underbrace{\Phi^T(s) \cdot w}_{= \hat{V}(s, w)}$$

$\left\{ \begin{array}{l} w - \text{parameter vector} \\ \Phi(s) - \text{the feature vector of } s \\ \hat{V}(s, w) \text{ is linear in } w \end{array} \right.$

广义线性模型

network

非线性系统

$V_\pi(s) \leq \hat{V}(s, w)$ called value approximation

表格中. $s_1 \rightarrow s_2$ s_2 in value 改变
其余不变

参数中. $s_1 \rightarrow s_2$ s_2 in value 改变
对应 $V(s, w)$ 中 w 改变
会导致 $s_1, s_2, V_{\text{value}}^{\text{PR}}$
变.

二. objective function

1. Optimal w so that $V(s, w)$ can best approximate $V_{\pi}(s)$ for every s .

2. Objective function:

$$J(w) = E[(V_{\pi}(s) - \hat{V}(s, w))^2]$$

① uniform distribution.

every state not equal important

② stationary distribution

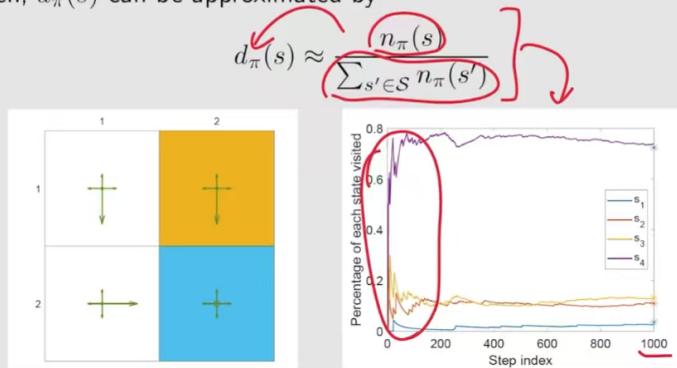
long run behavior $\{d_{\pi}(s)\}_{s \in S}$

$$J_{\pi} = E[(V_{\pi}(s) - \hat{V}(s, \pi))^2]$$

$$= \sum_{s \in S} d_{\pi}(s) (V_{\pi}(s) - \hat{V}(s, \pi))^2$$

Illustrative example:

- Given a policy shown in the figure.
- Let $n_{\pi}(s)$ denote the number of times that s has been visited in a very long episode generated by π .
- Then, $d_{\pi}(s)$ can be approximated by



Zhao

Figure: Long-run behavior of an ϵ -greedy policy with $\epsilon = 0.5$.

18

$$d_{\pi}^T = d_{\pi}^T P_{\pi}$$

$\xrightarrow{\text{状态转移矩阵}} \text{状态转移矩阵.}$

$$P(s'|s)$$

$$V_{\pi} = R_{\pi} + \gamma P_{\pi} V_{\pi}$$

根据 π

$$\Rightarrow P_{\pi} = \begin{bmatrix} 0.3 & 0.1 & 0.6 & 0 \\ 0.1 & 0.3 & 0 & 0.6 \\ 0.1 & 0 & 0.3 & 0.6 \\ 0 & 0.1 & 0.1 & 0.8 \end{bmatrix} \Rightarrow \text{代入 } d_{\pi}^T = d_{\pi}^T P_{\pi}$$

\downarrow

$$d_{\pi}^T = [0.045, 0.084, 0.133, 0.174]^T$$

3. Optimization algorithm:

gradient descent: $w_{k+1} = w_k - \alpha_k \nabla_w J(w_k)$

$$\begin{aligned}\nabla_w J(w_k) &= \nabla_w E[(V_\pi(s) - \hat{V}(s, w))^2] \\ &= E[\nabla_w (V_\pi(s) - \hat{V}(s, w))^2] \\ &= 2E[(V_\pi(s) - \hat{V}_\pi(s, w))(-\nabla_w \hat{V}(s, w))]\end{aligned}$$

Stochastic gradient \rightarrow true

$$w_{t+1} = w_t + \alpha_t (V_\pi(s_t) - \hat{V}(s_t, w_t)) \nabla_w \hat{V}(s_t, w_t)$$

replace $V_\pi(s_t)$

① Monte Carlo

$$w_{t+1} = w_t + \alpha_t (g_t - \hat{V}(s_t, w_t)) \nabla_w \hat{V}(s_t, w_t)$$

② TD learning:

$$w_{t+1} = w_t + \alpha_t [r_{t+1} + \gamma \hat{V}(s_{t+1}, w_t) - \hat{V}(s_t, w_t)] \nabla_w \hat{V}(s_t, w_t)$$

4. How to select function $\hat{V}(s, w)$

$$\Phi \quad \hat{V}(s, w) = \underline{\Phi}^T(s) w$$

② network

5. In the linear case:

$$\hat{V}(s, w) = \underline{\Phi}^T(s) w$$

We have $\nabla_w \hat{V}(s, w) = \underline{\Phi}(s)$

$$w_{t+1} = w_t + \alpha_t [r_{t+1} + \gamma \hat{V}(s_{t+1}, w_t) - \hat{V}(s_t, w_t)] \nabla_w \hat{V}(s_t, w_t)$$

$$w_{t+1} = w_t + \alpha_t [r_{t+1} + \gamma \underline{\Phi}^T(s_{t+1}) w_t - \underline{\Phi}^T(s_t) w_t] \underline{\Phi}(s_t)$$

TD-linear .

theoretical properties

6. Illustrative example. pg 12 .

III. Sarsa with function approximation.

$$W_{t+1} = W_t + \alpha_t [R_{t+1} + \gamma \hat{q}(S_{t+1}, a_{t+1}, W_t) - \hat{q}(S_t, a_t, W_t)] \nabla_w \hat{q}(S_t, a_t, W_t)$$

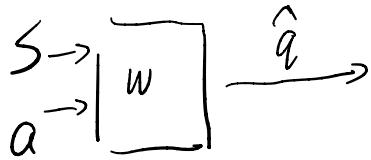
IV Q-learning with function approximation

$$W_{t+1} = W_t + \alpha_t [R_{t+1} + \gamma \max_{a \in A(S_{t+1})} \hat{q}(S_{t+1}, a, W_t) - \hat{q}(S_t, a_t, W_t)] \nabla_w \hat{q}(S_t, a_t, W_t)$$

3. Deep Q-learning

1. deep Q network.

$$W_{t+1} = W_t + \alpha_t [r_{t+1} + \gamma \max_{a \in A(S_{t+1})} \hat{q}(S_{t+1}, a, W_t) - \hat{q}(S_t, a_t, W_t)] \nabla_W \hat{q}(S_t, a_t, W_t)$$



$$J(w) = E \left[(R + \gamma \max_{a \in A(S')} \hat{q}(S', a, w) - \hat{q}(S, A, w))^2 \right]$$

2. 2 network.

① main network representing $\hat{q}(S, A, w)$

② target network $\hat{q}(S, A, w_T)$

$$J = E \left[(R + \gamma \max_{a \in A(S')} \hat{q}(S', a, w_T) - \hat{q}(S, A, w))^2 \right]$$

3. Experience. reply: 经验回放.

$$B \doteq \{(s, a, r, s')\}$$

混列 -> uniform distribution. 再拿出来训练

