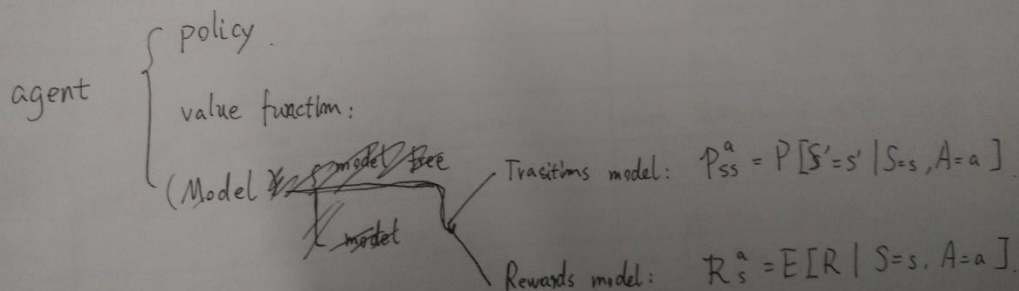
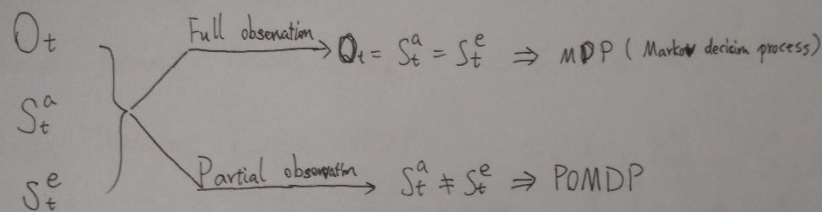


强化学习 第一讲



value based ~~agent~~ RL : { No Policy (或称为隐含)
Value function

policy based ~~agent~~ RL : { Policy
No value function

Actor Critic RL : { Policy
value function

Model free RL : { Policy and/or Value function
No Model

Model based RL : { Policy and/or Value function
Model

{ Exploration : 探索 (随机)
Exploitation : 开发 (最好的)

Prediction problem : evaluate the future.



Control problem : Find best policy

强化学习 第二讲

Markov Process (Markov Chain) : $\langle S, P \rangle$

Markov Reward Process (MRP) : $\langle S, P_{ss'}, R, \gamma \rangle$

Markov Decision Process (MDP) : $\langle S, A, P_{ss'}^a, R_s^a, \gamma \rangle$

	MRP $\langle S, P, R, \gamma \rangle$	MDP $\langle S, A, P, R, \gamma \rangle$
S 定义	finite set of state.	finite set of state.
A 定义	\emptyset	finite set of actions.

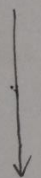
	MRP	MDP
P 定义	state transition probability matrix	state transition probability matrix.
P 方程	$P_{ss'} = P[S_{t+1} = s' S_t = s]$	$P_{ss'}^a = P[S_{t+1} = s' S_t = s, A_t = a]$
R 定义	reward function	reward function
R 方程	$R_s = E[R_{t+1} S_t = s]$	$R_s^a = E[R_{t+1} S_t = s, A_t = a]$
γ 定义	discount factor $\in [0, 1]$	discount factor $\in [0, 1]$
G_t 定义	$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$	$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$
$V(s)$	$V(s) = E[G_t S_t = s]$	$V_\pi(s) = E_\pi[G_t S_t = s]$
Bellman	$V(s) = E[G_t S_t = s] = E[R_{t+1} + \gamma V(S_{t+1}) S_t = s]$ $V(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} V(s')$	$V_\pi(s) = E_\pi[R_{t+1} + \gamma V_\pi(S_{t+1}) S_t = s]$ $V_\pi(s) = \sum_{a \in A} \pi(a s) q_\pi(s, a)$
$q_\pi(s, a)$	\emptyset	$q_\pi(s, a) = E_\pi[G_t S_t = s, A_t = a]$ $q_\pi(s, a) = E_\pi[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) S_t = s, A_t = a]$ $q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s')$
$V_*(s)$ 定义		$V_*(s) = \max_{\pi} V_\pi(s)$
$q_*(s, a)$ 定义		$q_*(s, a) = \max_{\pi} q_\pi(s, a)$
$V_*(s)$		$V_*(s) = \max_a q_*(s, a)$
$q_*(s, a)$		$q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s')$

强化学习 第二讲

Markov Process (Markov Chain): $\langle S, P \rangle$.



MRP: $\langle S, P_{ss'}, R, \gamma \rangle$



Bellman Equation: $\begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} V_1 \\ \vdots \\ V_n \end{bmatrix}$

iterative method $\left\{ \begin{array}{l} \text{Dynamic programming} \\ \text{Monte-Carlo evaluation} \\ \text{temporal-difference learning} \end{array} \right.$

MDP: $\langle S, A, P_{ss'}^a, R_s^a, \gamma \rangle$

Policy (π): $\pi(a|s) = P[A_t = a | S_t = s]$

State-value function: (following policy π)

$$V_{\pi}(s) = E_{\pi}[G_t | S_t = s] \xrightarrow{\text{Bellman}} E_{\pi}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s]$$

~~State~~ action-value function: (following policy π)

$$Q_{\pi}(s, a) = E_{\pi}[G_t | S_t = s, A_t = a] \xrightarrow{\text{Bellman}} E_{\pi}[R_{t+1} + \gamma Q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$

V和Q的关系 $\left\{ \begin{array}{l} V_{\pi}(s) = \sum_{a \in A} \pi(a|s) Q_{\pi}(s, a) \end{array} \right.$

$$Q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s')$$

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s') \right)$$

$$Q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \sum_{a' \in A} \pi(a'|s') Q_{\pi}(s', a')$$

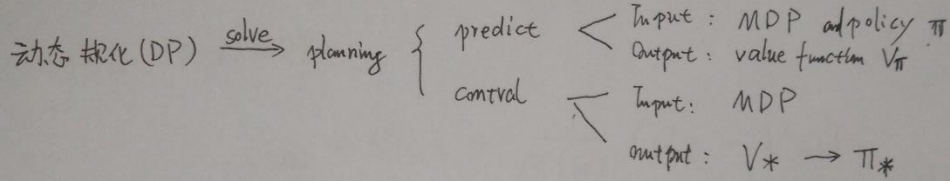
Optimal state-value function:

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

Optimal action-value function:

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

强化学习 第三讲 动态规划



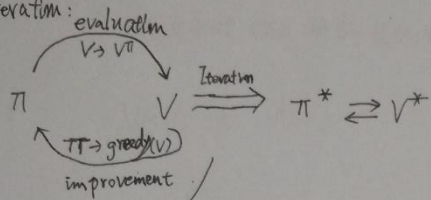
evaluation (using synchronous backups):

At each iteration $k+1$;

for all state $s \in S$:

update $V_{k+1}(s)$ from $V_k(s')$ where s' is a successor state of s .

Policy Iteration:



$$\text{greedy: } \pi'(s) = \underset{a \in A}{\operatorname{argmax}} q_\pi(s, a)$$

证明: (?)

收敛到 π^* 和 V^*

Value Iteration: At each iteration $k+1$:

For all state $s \in S$:

Update $V_{k+1}(s)$ from $V_k(s')$.

Problem	Bellman Equation	Algorithm
Predict	Bellman Expectation Equation	Iteration Policy Evaluation
Control	Bellman Expectation Equation + Greedy Policy Improvement	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration

强化学习第四讲 无模型预测

无模型 \Leftrightarrow 无 MDP

$$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (G_t - V(s_t))$$

\Downarrow 忘记过去很久的.

MC: $V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t))$

TD: $V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$
 \hookrightarrow Bellman Equation \rightarrow 在 MDP 中 TD 更有效.

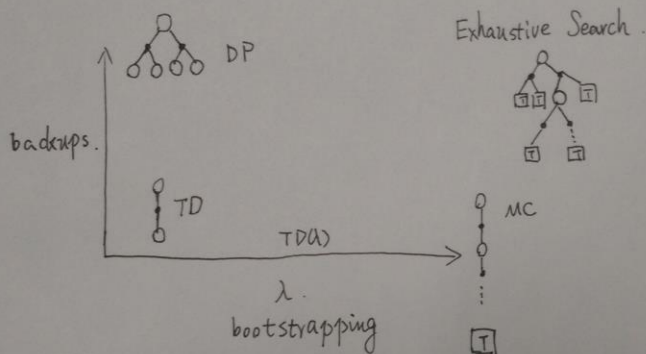
收敛于对数据解释最好的 MDP.

TD: $V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$

\hookrightarrow Bellman Equation \rightarrow 在 MDP 中 TD 更有效.

收敛于对数据解释最好的 MDP.

DP: $V(s_t) \leftarrow V(s_t) + E_{\pi} [R_{t+1} + \gamma V(s_{t+1})]$



TD(0). $G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t+n}^{(n)}$
 combine all n-step return

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t^\lambda - V(s_t))$$