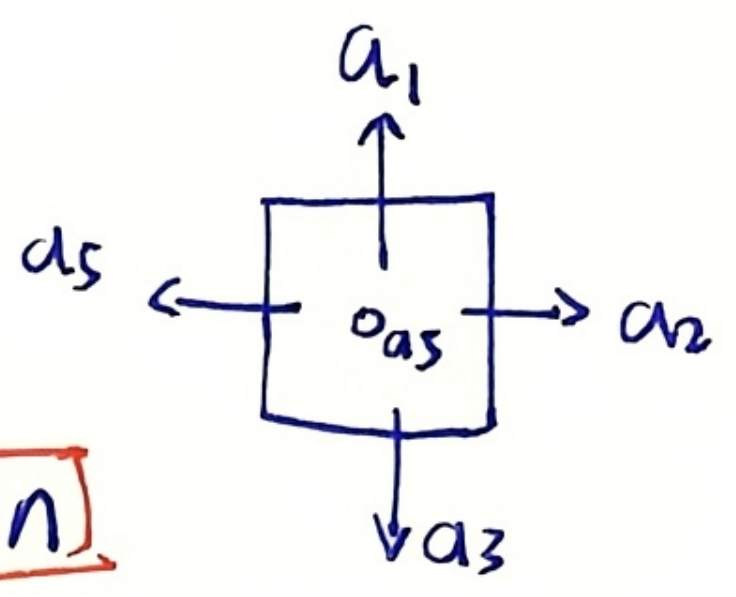


state space  $S = \{s_i\}_{i=1}^9$

action space  $A(s_i) = \{a_i\}_{i=1}^5$



$s_1$	$s_2$	$s_3$
$s_4$	$s_5$	$s_6$
$s_7$	$s_8$	$s_9$

$s_1 \xrightarrow{a_2} s_2$  : state transition

$s_1 \xrightarrow{a_1} s_1$  : 撞墙了.

→ defines the interaction with env.

two case {  $s_5 \xrightarrow{a_2} s_6$  accessible but with penalty  
 $s_5 \xrightarrow{a_2} s_5$  inaccessible

Tabular representation:

	$a_1$	$a_2$	$a_3$	...
$s_1$				
$s_2$				
⋮				

but only deterministic.

+ probability.

state transition probability.

$s_1 \xrightarrow{a_2} s_2$  :  $P(s_2 | s_1, a_2) = 1$

条件概率.

$P(s_i | s_1, a_2) = 0 \quad \forall i \neq 2$

policy.

what actions to take at a state

→	↓	←
→	↓	↓
→	→	○

different paths with different starting points

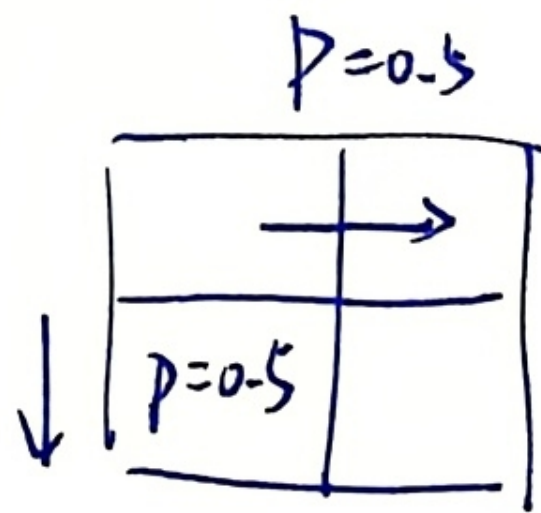
# Mathematical representation

$\pi \rightarrow$  for policy.

$$\left. \begin{array}{l} \pi(a_1 | s_1) = 0 \\ \pi(a_2 | s_1) = 1 \\ \pi(a_3 | s_1) = 0 \\ \dots \\ \dots \end{array} \right\} \text{plus} = 1.$$

1. deterministic

2. stochastic.



$$\left. \begin{array}{l} \pi(a_1 | s_1) = 0 \\ \pi(a_2 | s_1) = 0.5 \\ \pi(a_3 | s_1) = 0.5 \\ \dots \end{array} \right\} \text{plus} = 1.$$

random.rand(). 创造 0~1 随机数.

若  $x \leq 0.5$  take  $a_2$ .  $x > 0.5$  take  $a_3$ .

## ★ reward

human-machine interface

$> 0$  encouragement  
 $< 0$  punishment,  $>$  也可以正负互换.  
 $= 0$ . no punishment.

~~out of boundary. / forbidden cell = -1~~

$r_{\text{bound}} = -1$      $r_{\text{forbid}} = -1$      $r_{\text{target}} = +1$   
 otherwise  $r = 0$ .

## Mathematical description

$$\boxed{S_1 \xrightarrow{a_1}} \rightarrow r = -1$$

$$\begin{cases} P(r = -1 | s_1, a_1) = 1 \\ P(r \neq -1 | s_1, a_1) = 0 \end{cases}$$

reward transition. stochastic.  
正负也许可知. 但具体值未知

## Trajectory

扔轨迹  
Trajectory

a state-action-reward chain

$$S_1 \xrightarrow[r=0]{a_2} S_2 \xrightarrow[r=0]{a_2} S_5 \xrightarrow[r=0]{a_3} S_8 \xrightarrow[r=1]{a_2} S_9$$

return: the sum of them.  $0+0+0+1=1$

$$\text{return} = 1$$

return and reward: reward (R) 是 action 在 此刻 的反馈  
return (G) 是 reward 累加和

## Discounted return

引入 discount rate  $\gamma \in [0, 1)$

$$= 0 + \gamma^0 0 + \gamma^2 0 + \gamma^3 1 + \gamma^4 1 + \gamma^5 1 + \dots$$

$$= \gamma^3 (1 + \gamma + \gamma^2 + \dots) = \gamma^3 \frac{1}{1-\gamma} \leftarrow$$

可评估更近更近的 reward.  $\begin{cases} \gamma \text{ close close } 0 & \text{near future} \\ \gamma \text{ close } 1 & \text{far future.} \end{cases}$

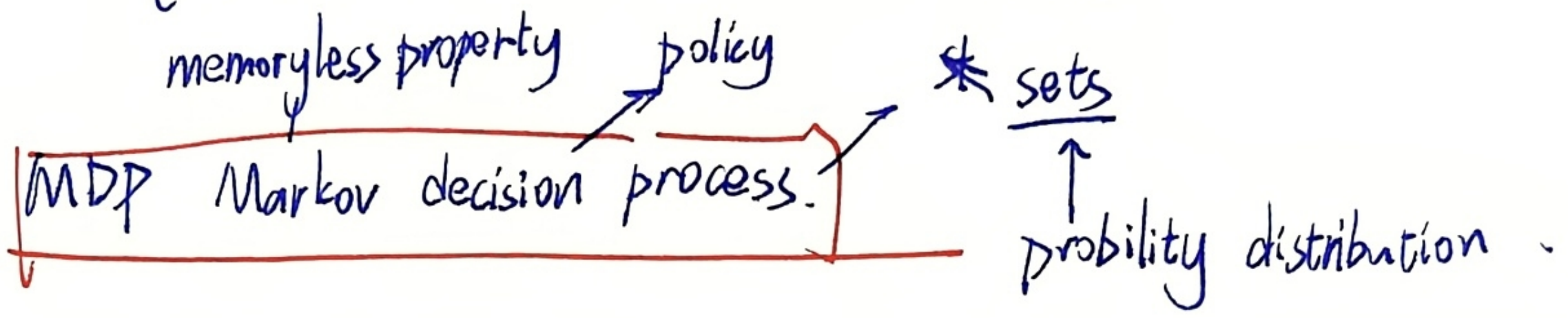
# Episode

有限步 finite 的 tasks.

continuing tasks vs Episode.

为了数学形式统一, 有以下两种方法统一:

- option 1. absorbing state. 仅有留在原地的 action, reward:
- option 2. normal state with policy. reward=1. 更一般化.



- sets:
    - state: the set of states  $S$
    - action: the set of actions  $A(s)$
    - reward: the set of rewards  $R(s, a)$
- MDP 无历史无记忆性的... 关于 policy... 以及 probability distribution 的状态转移 process

- probability:
  - state  $p(s'|s, a)$
  - reward  $p(r|s, a)$

- policy:
  - $s, a, \pi(a|s)$
  - $\sum_a \pi(a|s) = 1$

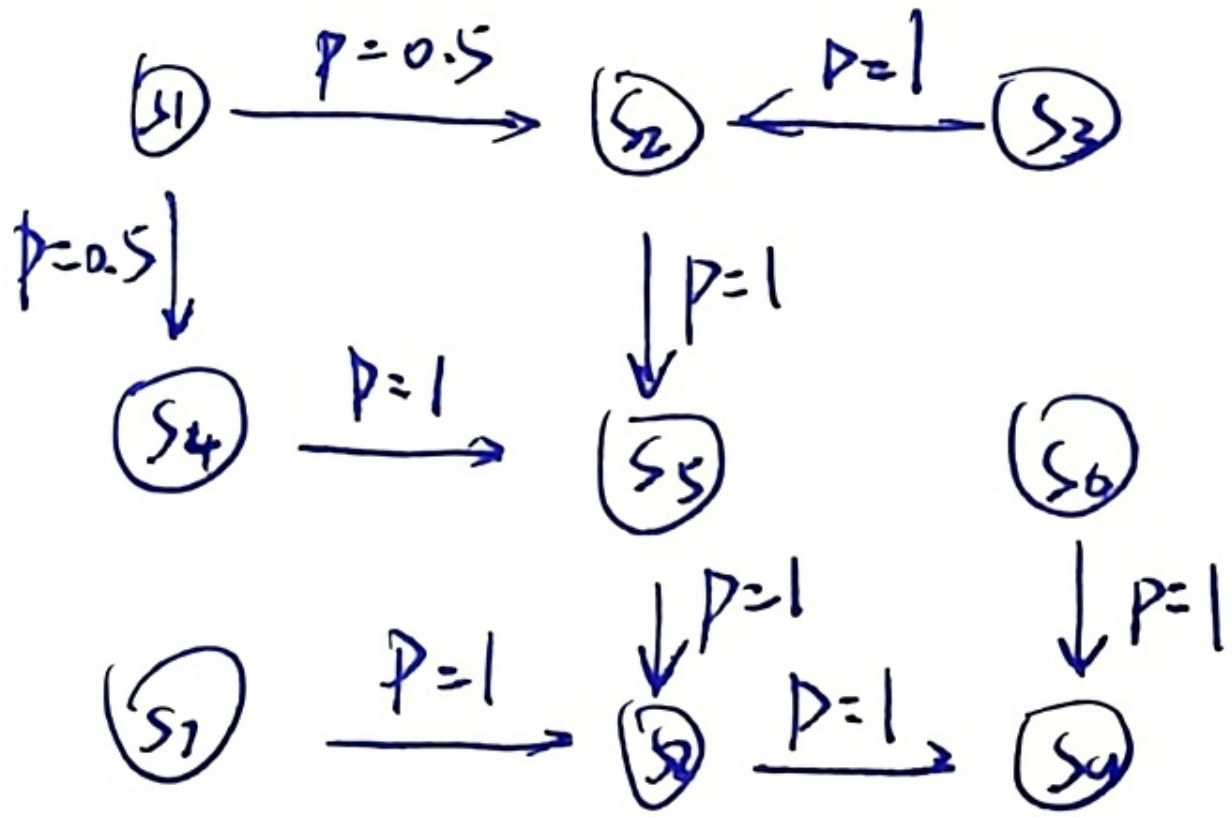
4. Markov property: memoryless property

~~$p(s_{t+1} | a_{t+1}, s_t)$~~

$p(s_{t+1} | \underline{a_{t+1}, s_t, \dots, a_1, s_0}) = p(s_{t+1} | \underline{a_{t+1}, s_t})$

$p(r_{t+1} | \underline{a_{t+1}, s_t, \dots, a_1, s_0}) = p(r_{t+1} | \underline{a_{t+1}, s_t})$

Markov process



Markov decision process becomes Markov process  
 once the policy is given!

↓  
 $\pi(a|s)$

MDD:  $P(s'|s, a)$

MP  $P(s'|s)$ .  $a \in \mathcal{A}(s)$

State value  ~~$V_{\pi}(s)$~~   $V_{\pi}(s)$ .

$$S_t \xrightarrow{A_t} R_{t+1}, S_{t+1}$$

$x_t, A_t, R_{t+1}$  are all random variables. 需引入  $E$

$$S_t \rightarrow A_t: \pi(A_t=a | S_t=s)$$

$$S_t, A_t \rightarrow R_{t+1}: p(R_{t+1}=r | S_t=s, A_t=a)$$

$$S_t, A_t \rightarrow S_{t+1}: p(S_{t+1}=s' | S_t=s, A_t=a)$$

↗  
转移概率.

trajectory:  $S_t \xrightarrow{A_t} R_{t+1}, S_{t+1} \xrightarrow{A_{t+1}} R_{t+2}, S_{t+2} \xrightarrow{A_{t+2}} R_{t+3} \dots$

discounted return:  $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$

↑  
random variables.

state-value ~~for~~  $V_{\pi}(s)$  =  $E[G_t | S_t=s]$ .

$$V(s, \pi) \quad V_{\pi}(s).$$

$V_{\pi}(s)$  vs  $G$ .

→  $V_{\pi}(s)$  是从某状态出发反获得所有可能回报的均值.

$\pi(a|s), p(r|s,a), p(s'|s,a)$  若都确定, 则  $V_{\pi}(s) = G$

Bellman Equation  $\left\{ \begin{array}{l} \text{state value} \\ \text{the Bellman equation} \end{array} \right.$

return 可评估策略好坏.

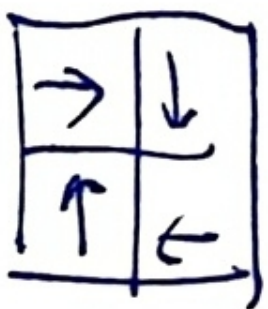
$$\begin{aligned} \text{return}_1 &= 0 + \gamma 1 + \gamma^2 1 + \dots \\ &= \frac{\gamma}{1-\gamma} \end{aligned}$$

$$\begin{aligned} \text{return}_2 &= -1 + \gamma 1 + \gamma^2 1 + \dots \\ &= -1 + \frac{\gamma}{1-\gamma} \end{aligned}$$

$$\begin{aligned} \text{return}_3 &= 0.5 \left( -1 + \frac{\gamma}{1-\gamma} \right) + 0.5 \left( \frac{\gamma}{1-\gamma} \right) \\ &= -0.5 + \frac{\gamma}{1-\gamma} \end{aligned}$$

$\text{return}_1 > \text{return}_3 > \text{return}_2$ .

Method 1:



$$V_1 = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$$

$$V_2 = r_2 + \gamma r_3 + \gamma^2 r_4 + \dots$$

$$V_3 = r_3 + \gamma r_4 + \gamma^2 r_1 + \dots$$

$$V_4 = r_4 + \gamma r_1 + \gamma^2 r_2 + \dots$$

The returns  
rely on each other.

$$V_1 = r_1 + \gamma (r_2 + \gamma r_3 + \dots) = r_1 + \gamma V_2$$

$$V_2 = r_2 + \gamma (r_3 + \gamma r_4 + \dots) = r_2 + \gamma V_3$$

$$V_3 = r_3 + \gamma (r_4 + \gamma r_1 + \dots) = r_3 + \gamma V_4$$

$$V_4 = r_4 + \gamma (r_1 + \gamma r_2 + \dots) = r_4 + \gamma V_1$$

Bootstrapping 自举

$$\begin{aligned}
V_{\pi}(s) &= E[G_t | S_t = s] \\
&= E[R_{t+1} + \gamma G_{t+1} | S_t = s] \\
&= \underbrace{E[R_{t+1} | S_t = s]}_{(1)} + \gamma \underbrace{E[G_{t+1} | S_t = s]}_{(2)}
\end{aligned}$$

(1):  $E[R_{t+1} | S_t = s] = \sum_a \pi(a|s) E[R_{t+1} | S_t = s, A_t = a]$   
immediate rewards  $= \sum_a \pi(a|s) \sum_r p(r|s, a) r$

(2):  $E[G_{t+1} | S_t = s] = \sum_{s'} E[G_{t+1} | S_t = s, S_{t+1} = s'] p(s'|s)$   
(马尔可夫性质)  $= \sum_{s'} E[G_{t+1} | S_{t+1} = s'] p(s'|s)$   
 $= \sum_{s'} V_{\pi}(s') p(s'|s)$   
 $= \sum_{s'} V_{\pi}(s') \sum_a p(s'|s, a) \pi(a|s)$

对所有状态成立: Bellman:

$$V_{\pi}(s) = \sum_a \pi(a|s) \left[ \sum_r p(r|s, a) r + \gamma \sum_{s'} p(s'|s, a) V_{\pi}(s') \right]$$

$\Delta$   $\Delta$   
 $s$ -state  $v$  immediate future  $s'$ -state  $v$

Bootstrapping 记算  $\pi(a|s)$  是给出的. policy evaluation  
 $p(r|s, a)$   $p(s'|s, a)$ . dynamic model 环境.

# 贝尔曼方程

不同策略的条件概率累加 (加权).

[前不同策略下的 reward 的累加. 加权].

+ [折扣后未来状态的状态函数的累加. 加权].

累加权与实际项.

$$\# V_{\pi}(s) = \sum_a \pi(a|s) \left[ \underbrace{\sum_r P(r|s,a)}_{\text{大累加套}} r + \gamma \underbrace{\sum_{s'} P(s'|s,a)}_{\text{大累加套}} V_{\pi}(s') \right]$$

本质是  $V_{\pi}(s) = r + \gamma V_{\pi}(s')$  的一般化.

由  $V_{\pi}(s) = r + \gamma r_1 + \gamma^2 r_2 + \gamma^3 r_3 \dots$  得到.

$$\sum_a \pi(a|s).$$

$$\sum_r P(r|s,a).$$

引入不同 action/策略. (不同策略下) 引入不同的 reward 概率. 与  $r$ .

$$\sum_{s'} P(s'|s,a)$$

(同策略下). 去到  $s'$  新状态的概率. 与  $V_{\pi}(s')$

# policy evaluation.

a policy. finding out the corresponding state values is called policy evaluation.

## solve Bellman

迭代\*

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

closed-form:  $V_{\pi} = (I - \gamma P_{\pi})^{-1} r_{\pi}$

iterative solution:  $V_{k+1} = r_{\pi} + \gamma P_{\pi} V_k$

$$V_0 \rightarrow V_1$$

$$V_1 \rightarrow V_2$$

$$V_2 \rightarrow V_3$$

sequence  $\{V_0, V_1, V_2, \dots\}$

$$V_k \rightarrow V_{\pi} = (I - \gamma P_{\pi})^{-1} r_{\pi}, \quad k \rightarrow \infty$$

$V_{\pi} = r_{\pi} + \gamma P_{\pi} V_{\pi}$ . 设  $F(V_{\pi}) = r_{\pi} + \gamma P_{\pi} V_{\pi}$

则  $V_{\pi}$  为  $F(V_{\pi})$  的不动点.

$$V_{k+1} = r_{\pi} + \gamma P_{\pi} V_k$$

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

$$|V_{k+1} - V_{\pi}| = \gamma P_{\pi} |V_k - V_{\pi}|$$

新误差  $\leq \gamma \times$  旧误差.



策略只 take  $a_2$  情况下. 其他  $a$  也都是可以计算的. 不等于 0.

① 先算 state values 后算 action values.

②\* 直接算 action values  $\begin{cases} \text{with models} \\ \text{without models.} \end{cases}$

**Bellman Optimality Equation** BOE.

$$a^* = \operatorname{argmax}_a q_{\pi}(s, a).$$

$$\pi_{\text{new}}(a|s) = \begin{cases} 1, & a = a^* \\ 0, & a \neq a^* \end{cases}$$

$\pi^*$   
贪心策略.

greedy  $\pi^*$ .

定义:  $V_{\pi_1}(s) \geq V_{\pi_2}(s)$ . for all  $s \in S$ .

$\pi_1$  is better than  $\pi_2$ .

"A policy  $\pi^*$  is optimal if  $V_{\pi^*}(s) \geq V_{\pi}(s)$  for all  $s$  and for any other policy  $\pi$ "

**BOE**

$$V(s) = \max_{\pi} \sum_a \pi(a|s) \left( \sum_r P(r|s, a) r + \gamma \sum_{s'} P(s'|s, a) V(s') \right)$$

$$= \max_{\pi} \sum_a \pi(a|s) q(s, a).$$

$$V = \max_{\pi} (r_{\pi} + \gamma P_{\pi} V).$$

(BOE).

$$V = \max_{\pi} (r\pi + \gamma R)$$

$$V = \max_{\pi} (r\pi + \gamma P_{\pi} V). \quad \text{tricky and elegant.}$$

How to solve  $\max_{\pi} \sum_a \pi(a|s) q(s, a)$

对于  $q$ . (action value) 最大的. 将策略  $\pi(a|s)$  都向大的  $q$  分配. 可以确保  $V$  的 optimal.

$$\begin{aligned} V &= \max_{\pi} \sum_a \pi(a|s) \left( \sum_r p(r|s, a) r + \gamma \sum_{s'} p(s'|s, a) V(s') \right) \\ &= \max_{\pi} \sum_a \pi(a|s) q(s, a). \end{aligned}$$

considering  $\sum_a \pi(a|s) = 1$ .  $\max_{\pi} \sum_a \pi(a|s) q(s, a) = \max_{a \in A_s} q(s, a)$

取 optimality when  $\pi(a|s) = \begin{cases} 1, & a = a^* \\ 0, & a \neq a^* \end{cases}$

$$a^* = \operatorname{argmax}_a q(s, a)$$

---

3.3 How to solve BOE.

solve BOE:

$$f(v) = \max_{\pi} (r_{\pi} + \gamma P_{\pi} v) \quad v = f(v) \quad *$$

$$[f(v)]_s = \max_{\pi} \sum_a \bar{\pi}(a|s) q(s, a)$$

收缩  
Contraction mapping theorem.

fixed point:  $x \in X$ .  $f(x) = x$ . 即 fixed point.

$$\|f(x_1) - f(x_2)\| \leq \gamma \|x_1 - x_2\|$$

① 一定有  $x^*$  使  $f(x^*) = x^*$

② unique.

③  $x_{k+1} = f(x_k)$   $x_k \rightarrow x^*$  as  $k \rightarrow \infty$ .

Convergence rate is exponentially fast

对矩阵同理.

BOE  $v = f(v) = \max_{\pi} (r_{\pi} + \gamma P_{\pi} v)$ .

是 contraction mapping.

# Policy optimality

$V^*$  is optimality equation to BOE.

$$V^* = \max_{\pi} (r_{\pi} + \gamma P_{\pi} V^*).$$

suppose  $\pi^* = \arg \max_{\pi} (r_{\pi} + \gamma P_{\pi} V^*)$

$$\therefore V^* = r_{\pi^*} + \gamma P_{\pi^*} V^*.$$

$V^* = V_{\pi^*}$ . 是特殊的 BE.

$V^* \geq V_{\pi}$ .  $\pi$  为 optimality policy.

$$\pi^* \quad \pi^*(a|s) = \begin{cases} 1 & \leftarrow a = a^*(s) \\ 0 & \uparrow a \neq a^*(s). \end{cases} \quad \pi^* \text{ deterministic greedy.}$$

## BOE

Reward design:  $r$

System model:  $p(s'|s, a), p(r|s, a)$

Discount rate:  $\gamma$ .

$\vec{r} \rightarrow a r + b$ .  $R$  要  $a > 0$ .

$$V' = a V + \frac{b}{1-\gamma}.$$

state value 唯一, 但  $\pi$  不一定唯一.

# Value 值迭代算法

(比较蒙. 例子多看几遍)

step 1. policy update.

step 2. value update.  $V_k$  并非 state value.

$$V_k(s) \rightarrow q_k(s, a) \rightarrow \text{greedy policy } \pi_{k+1}(a|s) \rightarrow \text{new value } V_{k+1} \\ = \max_a q_k(s, a).$$

老师的例子讲得很好.

有 Q-table. \* 对应的  $q$  由  $a, v$  决定.

$v$  随机一下. 可以找到  $q$  大的  $v_1, v_2, v_3, \dots$

~~将  $v_1$~~ , 以及  $v_1, v_2, v_3$  所 take 的 action

↓  
value update.

↓  
policy update.

↓ 对于每一个 state 都有当前最优 value.

进入下一个迭代

继续又能算出 Q-table 中每个 state 对应的最优 Q

又能更新 value 与 policy.