

action: 动作
 state: 状态.
 policy: 策略函数

$$\pi(s, a)$$

$\pi(a|s) = P(A=a | S=s)$ 即在当前状态 s 下选取动作 a 的概率
 随机的原因: 如果在 s 下 a 是确定的, 那么便是熟练工, 很容易让对手赢

Reward: 奖励:

old state $\xrightarrow{\text{action}}$ new action.

$P(s'=s' | s=s, A=a)$ = 在当前状态 s 下采取行动 a 产生新状态 s' 的概率

① 动作是随机的 ② 状态转移是随机的

状态 s , 动作 a , 奖励 r, \dots 称为轨迹

Return = 从 t 时刻开始到结束时的奖励之和.

$$U_t = R_t + R_{t+1} + \dots$$

discounted return = 折扣率回报:

$$U_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} \dots$$

动作价值函数 (策略 π 的)

$$Q_\pi(s_t, a_t) = E[U_t | S_t = s_t, A_t = a_t]$$

因为 U_t 依赖于动作 $A_t, A_{t+1}, A_{t+2}, \dots$ 和状态 S_t, S_{t+1}, \dots

最优动作价值函数:

$$Q^*(s_t, a_t) = \max_{\pi} Q_\pi(s_t, a_t)$$

状态价值函数:

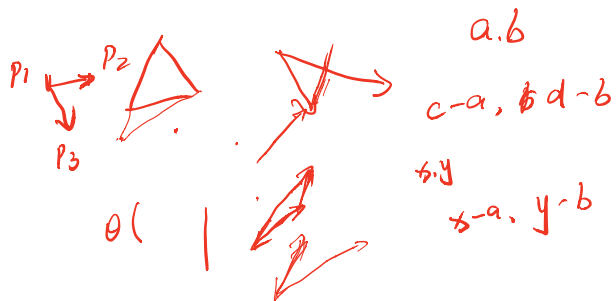
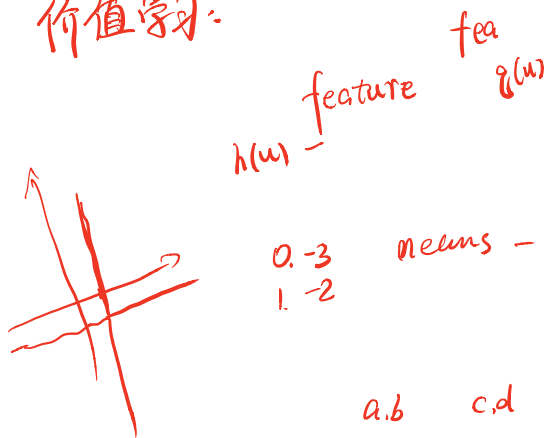
$$V_\pi(s_t) = E_A[Q_\pi(s_t, A)] = \sum_a \pi(a|s_t) Q_\pi(s_t, a)$$

其中 $A \sim \pi(\cdot | s_t)$

$$= \int \pi(a|s_t) Q_\pi(s_t, a) da$$

↓
看当前的状态好不好

价值学习:



$$(c-a)(x-a) + (d-b)(y-b)$$