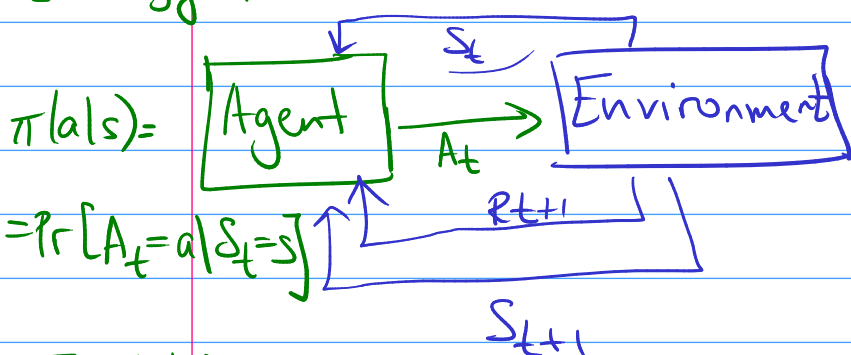


① MDP - Markov decision process

$$\mathcal{S}, A(s), r \in \mathbb{R}$$

strategy π



$$\pi(a|s) = \Pr[A_t = a | S_t = s]$$

$$\sum_a \pi(a|s) = 1 \quad \forall s$$

$$p(s'|s,a) = \sum_z p(s',z|s,a)$$

$$r(s,a) = \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a] = \sum_z r \cdot \underbrace{\sum_{s'} p(s',z|s,a)}_{p(z|s,a)}$$

$$\overbrace{S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots}^{\text{dynamics}}$$

$$p(s',z|s,a) =$$

$$= \Pr[S_{t+1} = s', R_{t+1} = z | S_t = s, A_t = a]$$

$$\sum p(s',z|s,a) = 1 \quad \forall s,a$$

$|S|=1$ - bandit

$|S|=1$ + korrekte c

$r(a/c)$

② Returns

$t: R_{t+1}, R_{t+2}, R_{t+3}, \dots$

return

discount

Episodic - T

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

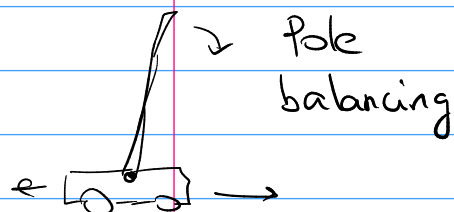
Continuous

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^k R_{t+k+1} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

M.S. $T=\infty$

M.S. $\gamma=1$



+1 t $\gamma=1$

= 1:gamma

$\gamma < 1$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots)$$

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

③ Value function

Agent

$$\pi(a|s) = \Pr[A_t = a | S_t = s]$$

$$\begin{aligned} V_{\pi}(s) &= E_{\pi}[G_t | S_t = s] = \\ &= E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right] \\ Q_{\pi}(s, a) &= E_{\pi}[G_t | S_t = s, A_t = a] = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a\right] \end{aligned}$$

$$V_{\pi}(s) = \sum_a \pi(a|s) \cdot Q_{\pi}(s, a)$$

$$\begin{aligned} Q_{\pi}(s, a) &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] = \\ &= \underbrace{E_{\pi}[R_{t+1} | S_t = s, A_t = a]}_{r(s, a)} + \gamma E_{\pi}[G_{t+1} | S_t = s, A_t = a] \end{aligned}$$

$$\begin{aligned} &\text{where } S_{t+1} = s' \\ &p(s' | s, a) \end{aligned}$$

$$Q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s'} p(s' | s, a) \cdot V_{\pi}(s')$$

④

Bellman equations

$$V_{\pi}(s) = E_{\pi}[G_t | S_t = s] =$$

$$= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] =$$

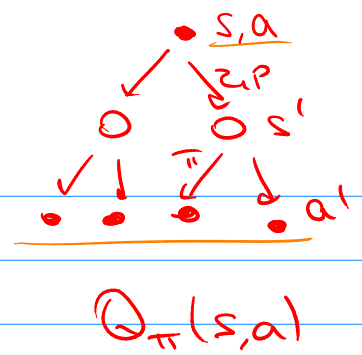
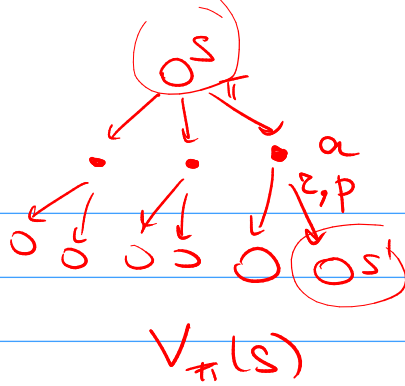
$$= \sum_a \pi(a|s) \cdot \sum_{s'} \sum_z p(s', z | s, a) [z + \gamma E_{\pi}[G_{t+1} | S_{t+1} = s']]$$

$$\forall s \quad V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} \sum_z p(s', z | s, a) [z + \gamma V_{\pi}(s')]$$

Bellman equations

$$\forall s, a \quad Q_{\pi}(s, a) = \sum_{s'} \sum_z p(s', z | s, a) [z + \gamma \sum_{a'} \pi(a' | s') Q_{\pi}(s', a')]$$

backup diagram

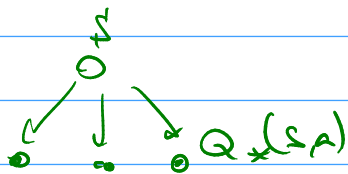


⑤ Optimal strategies

π_* - optimal policy, i.e. $\forall \pi \forall s \quad \boxed{V_{\pi_*}(s)} \geq V_{\pi}(s)$

$\pi: s \rightarrow \pi(a|s)$

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



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

$$\pi_*(s) = \operatorname{argmax}_a \boxed{Q_*(s,a)}$$

$$\underline{V_*(s)} = \max_{\pi} V_{\pi}(s) = \max_{\pi} \left[\sum_a \pi(a|s) \left[\sum_{s',z} p(s',z|s,a) [-] \right] \right] =$$

$$a = \operatorname{argmax}_a \max_{\pi} \sum_{s',z} p(-) [r + \gamma V_{\pi}(s')] =$$

$$= \operatorname{argmax}_a \sum_{s',z} p(s',z|s,a) [r + \gamma V_*(s')]$$

$\forall s$

$$V_*(s) = \max_{\pi} V_{\pi}(s) = \max_a \left(\sum_{s',z} p(s',z|s,a) [r + \gamma V_*(s')] \right) \quad \left. \begin{array}{l} \text{Bellman} \\ \text{equation} \end{array} \right\}$$

$$Q_*(s,a) = \max_{\pi} Q_{\pi}(s,a) = \sum_{s',z} p(s',z|s,a) [r + \gamma \max_{\pi} V_{\pi}(s')]$$

$$\max_{a'} \max_{\pi} Q_{\pi}(s',a')$$

$$Q_*(s,a) = \sum_{s',z} p(s',z|s,a) [r + \gamma \max_{a'} Q_*(s',a')]$$

$$\bar{x} = f(\bar{x})$$

$$\bar{x}^0, \bar{x}^1 = f(\bar{x}^0), \bar{x}^2 = f(\bar{x}^1), \dots$$

$$\bar{x}^{k+1} = f(\bar{x}^k)$$

$$\bar{x} = A\bar{x} - (I-A)\bar{x}$$

Policy evaluation

π :

$$V_{\pi}^{(k+1)}(s) := \text{Bellman}(V_{\pi}^{(k)}(s))$$

π_x

$$V_x^{(k+1)}(s) \rightarrow \dots$$

max

⑥ Policy improvement

Impr.

$$\pi \rightarrow V_{\pi}(s), Q_{\pi}(s,a)$$

Policy impr. theorem

Hypos π, π' - greedy w.r.t. Q_{π}

True $\forall s$

$$\forall s \quad Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \Rightarrow V_{\pi'}(s) \geq V_{\pi}(s)$$

1-bd: $V_{\pi}(s) \leq Q_{\pi}(s, \pi'(s)) =$

$$= E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = \pi'(s)] =$$

$$= E_{\pi'}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s] \leq$$

$$\leq E_{\pi'}[R_{t+1} + \gamma Q_{\pi}(S_{t+1}, \pi'(S_{t+1})) | S_t = s] = \dots$$

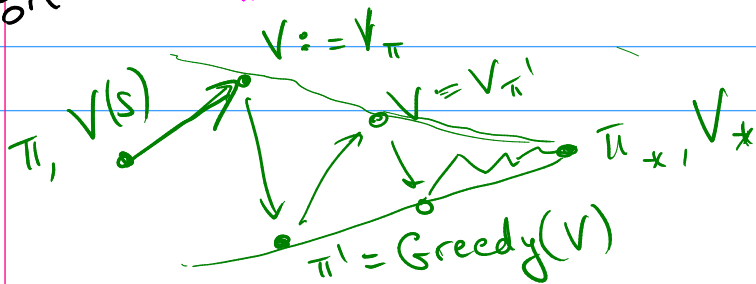
$$\leq \dots E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 Q_{\pi}(S_{t+2}, \pi'(S_{t+2})) | S_t = s]$$

$$\leq \dots E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s] = V_{\pi'}(s)$$

$$\pi'(s) = \operatorname{argmax}_a Q_{\pi}(s, a) = \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

Policy iteration

$$Q_{\pi}(s, \pi'(s)) \geq Q_{\pi}(s, \pi(s)) = V_{\pi}(s)$$



$$\pi' = \pi \Rightarrow \pi' = \pi = \pi_*$$

$$\pi(s) = \operatorname{argmax}_a Q_{\pi}(s, a)$$

DP
dynamic
programming

Value
iteration

$$V^{(k+1)}(s) = \max_a \sum_{s', z} p(s', z | s, a) [z + \gamma \cdot V^{(k)}(s')]$$

$V_{\pi}(s)$ $p(\dots)$

- $V(s) \leftarrow \text{random}$

- loop:

- Unpack $\pi, S_0, A_0, R_1, S_1, A_1, \dots, R_T$

- $G \leftarrow 0$

- for $t = T-1, T-2, \dots, 0$:

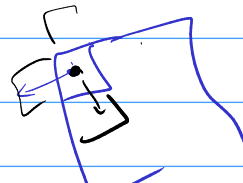
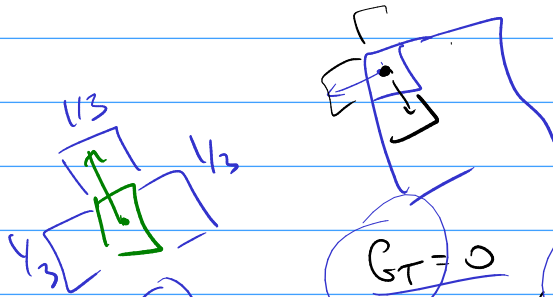
- $G \leftarrow \gamma G + R_{t+1}$

- compute G & enqueue $\text{Returns}(S_t)$

- $V(S_t) := \text{Avg}(\text{Returns}(S_t))$

- $V_{\pi}(s) := V(s)$

$\rightarrow \pi$



$G_T = 0$

$G_{T-1} = R_T$

$G_{T-2} = R_{T-1} + \gamma R_T$



first-visit MC
every-visit MC

$G = R_1 + \gamma R_2 + \gamma^2 R_3 + \dots + \gamma^{T-1} R_T$