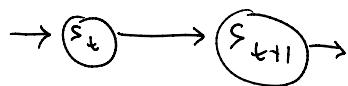


$$\mathcal{S} = \{1, 2\} \leftarrow$$

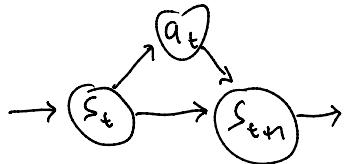
$$\text{MP: } P(S_{t+1}=j | S_t=i) := A_{ij}$$

MLE

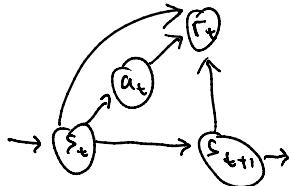
$$\hat{A}_{ij} = \frac{n_{ij}}{\sum_j n_{ij}}$$



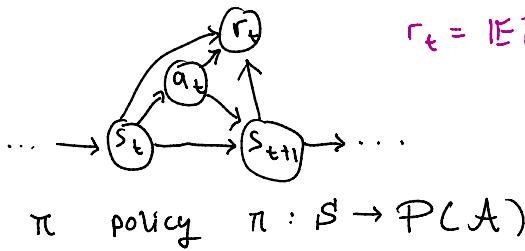
MDP:



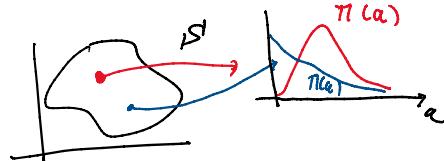
$$P(S_{t+1}=j | S_t=i, a_t=a) := A_{ij}^a$$



$$P(S_{t+1}=j, r_t=r | a_t=a, S_t=i) := A_{ij}^a r$$



$$r_t = \mathbb{E}[R_t | S_t=s]$$



$$\rightarrow V^\pi(s) := \mathbb{E}[r_1 + \gamma r_2 + \dots | S_0=s], \quad \gamma \in (0, 1)$$

state value of the policy \$\pi\$, being in state \$s\$

$$Q^\pi(s, a) := \mathbb{E}[r_1 + \gamma r_2 + \dots | S_0=s, a_0=a]$$

state action value of the policy \$\pi\$, being in state \$s\$ & taking action \$a\$.

Bellman principle

$$\begin{aligned} V^\pi(s) &:= \mathbb{E}[r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | S_0=s] \\ &= \mathbb{E}[r_1 + \underbrace{\mathbb{E}[r_2 + \gamma r_3 + \gamma^2 r_4 + \dots | S_1]}_{\gamma \mathbb{E}[r_2 + \gamma r_3 + \gamma^2 r_4 + \dots | S_1]} | S_0=s] \end{aligned}$$

$$\gamma \underbrace{[E[r_2 + \gamma' r_3 + \gamma^2 r_4 + \dots | s_1]]}_{V^\pi(s_1)}$$

$$V^\pi(s) = E[r_1 + \gamma V^\pi(s_1) | s_0 = s]$$

Bellman equation

$$s^1 = \{1, 2, \dots, T\}$$

$$\vec{V}^\pi = \begin{pmatrix} V^\pi(1) \\ V^\pi(2) \\ \vdots \\ V^\pi(T) \end{pmatrix} \quad \begin{aligned} & E[r_1 | s_0 = s] \\ \Rightarrow & =: \begin{pmatrix} E[r_1 | s_0 = 1] \\ E[r_1 | s_0 = 2] \\ \vdots \\ E[r_1 | s_0 = T] \end{pmatrix} = \int r(a, s) \pi(a | s) da \end{aligned}$$

$$\begin{aligned} \vec{V}_s^\pi &= V^\pi(s) = \sum_{s', a} (r(a, s) + \gamma V^\pi(s')) \underbrace{\Pr(a, s' | s)}_{\Pr(s' | a, s) \Pr(a | s)} da ds' \\ &= \sum_{s', a} (r(a, s) + \gamma V^\pi(s')) \underbrace{A_{s, s'}^a}_{\Pr(s' | a, s) / \Pr(a | s)} \underbrace{\pi(a | s)}_{\Pr(a | s)} \\ &= \vec{r}_s + \gamma \sum_{s'} A_{s, s'} \vec{V}_{s'}^\pi \end{aligned}$$

$$\vec{V}^\pi = \vec{r} + \gamma A \vec{V}^\pi$$

$$\Rightarrow (1 - \gamma A) \vec{V}^\pi = \vec{r}$$

$$\Rightarrow \boxed{\vec{V}^\pi = (1 - \gamma A)^{-1} \vec{r}}$$

analogously we have

$$\boxed{Q^\pi(s, a) = E[r_1 + \gamma Q^\pi(s_1, a_1) | s_0 = s, a_0 = a]}$$

$$Q^\pi(s, a) = \mathbb{E} [r_1 + \gamma Q^\pi(s_1, a_1) \mid s_0 = s, a_0 = a]$$

$$V^\pi(s) = \mathbb{E}_{\pi}^{\mathbb{P}_\pi} [r(s, a) + \gamma V^\pi(s') \mid s = s]$$

$$V^\pi(s) = \mathbb{E}_{\pi}^{\mathbb{P}_\pi} [Q^\pi(s, a) \mid s = s]$$

$$\mathbb{E}_{\pi}^{\mathbb{P}_\pi} [r(a, s) + \gamma r(a', s') + \dots \mid s = s, a = a]$$

Suppose I take an arbitrary action at $t=0$,

then follow the optimal policy π^*

call this $\tilde{\pi}$

$$V^{\tilde{\pi}}(s) = \mathbb{E}_{\tilde{\pi}}^{\mathbb{P}_{\tilde{\pi}}} [r(s, a) + \gamma \underbrace{V^{\tilde{\pi}}(s')}_{= V^{\pi^*}(s')} \mid s = s]$$

$$V^{\tilde{\pi}}(s) = \mathbb{E}_{\tilde{\pi}}^{\mathbb{P}_{\tilde{\pi}}} [r(s, a) + \gamma V^{\pi^*}(s') \mid s = s]$$

$$V^{\pi^*}(s) = \max_{\tilde{\pi} \in \mathcal{P}} \mathbb{E}_{\tilde{\pi}}^{\mathbb{P}_{\tilde{\pi}}} [r(s, a) + \gamma V^{\pi^*}(s') \mid s = s]$$

Bellman optimality equation

$$\begin{aligned}
 Q^{\pi^*}(s, a) &= \mathbb{E}^{\mathbb{P}_{\pi^*}} \left[r(s, a) + \gamma V^{\pi^*}(s') \mid s=s, a=a \right] \\
 &= \mathbb{E}^{\mathbb{P}_{\pi^*}} \left[r(s, a) + \gamma \mathbb{E}^{\mathbb{P}_{\pi^*}} \left[Q^{\pi^*}(s', a') \mid s' \right] \mid s=s, a=a \right]
 \end{aligned}$$

Suppose $\pi(a|s)$ are "deterministic" policies

$$\begin{aligned}
 V^{\pi^*}(s) &= \max_{a \in A} \mathbb{E}^{\mathbb{P}_{\pi^*}} \left[r(s, a) + \gamma V^{\pi^*}(s') \mid s=s \right] \\
 Q^{\pi^*}(s, a) &= \mathbb{E}^{\mathbb{P}_{\pi^*}} \left[r(s, a) + \gamma \max_{a' \in A} Q(s', a') \mid s=s, a=a \right]
 \end{aligned}$$

$$V^{\pi^*}(s') = \max_{a' \in A} Q^{\pi^*}(s', a')$$

$$\begin{aligned}
 \pi \rightarrow V^\pi(s) \xrightarrow[\text{greedy}]{} \pi' &= \arg \max_{a \in A} \mathbb{E}^{\mathbb{P}_{\pi}} [\dots] \\
 &\rightarrow V^{\pi'}(s)
 \end{aligned}$$

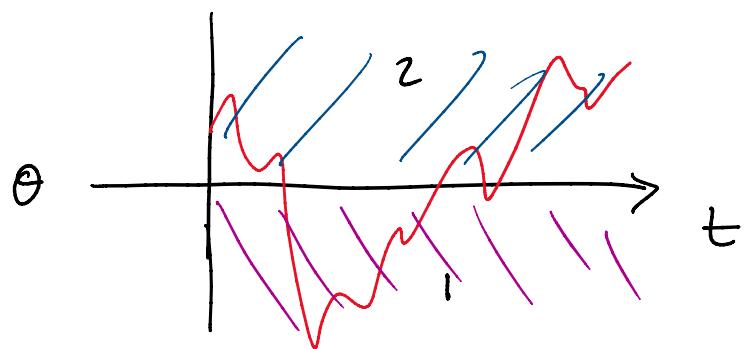
$$\Rightarrow V^{\pi'}(s) \geq V^\pi(s)$$

Define the Bellman operator $T^\pi: F_S \mapsto F_S$

$$f \in F_S \quad \begin{matrix} \nearrow \text{for deterministic policies} \\ a(s) \end{matrix} \\
 (T^\pi f)(s) = \mathbb{E}^{\mathbb{P}_{\pi}} \left[r(s, a) + \gamma f(s') \mid s=s \right]$$

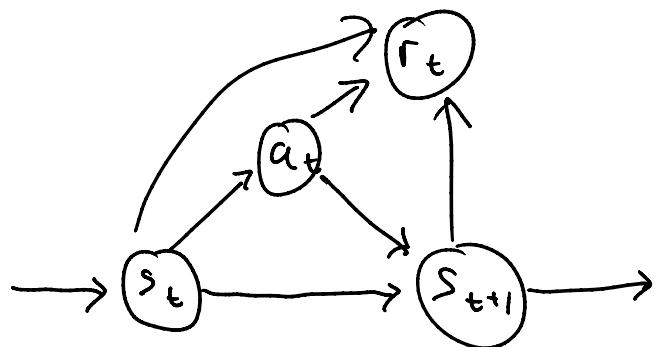
$$\text{claim: } \lim (T^\pi)^n f = V^\pi$$

$$\text{claim : } \lim_{n \rightarrow \infty} (T^n f) = V^n$$



$$A = \begin{pmatrix} 0.4 & 0.6 \\ 0.6 & 0.4 \end{pmatrix}$$

$A = \{$ do nothing, 1
long position, 2
short position $\} \{ 3$



rewards

initial policy

		a		
		1	2	3
s'		1	0	0
1	1	0	0	0
1	2	0	+1	-1
2	1	0	-1	+1
2	2	0	0	0

		π
		1
1	1	1
2	1	1

$$V^\pi(s) = \mathbb{E}^{\mathbb{P}_\pi} [r(s, a) + \gamma V^\pi(s') \mid s=s]$$

$$\begin{aligned} V^{\pi^{(k+1)}}(s) &= (T^\pi V^{\pi^{(k)}})(s) \\ &= \mathbb{E}^{\mathbb{P}_\pi} [\underbrace{r(s, a)}_{\uparrow} + \underbrace{\gamma V^{\pi^{(k)}}(s')}_{\uparrow} \mid s=s] \end{aligned}$$

sample from \mathbb{P}_π to get $(s^{(t)}, a^{(t)}, r^{(t)}, s'^{(t)})$

$$V_{t+1}^\pi(s_t) = \mathbb{E}^{\mathbb{P}_\pi} [r(s_t, a_t) + \gamma V_t^\pi(s_{t+1}) \mid s_t]$$

Temporal - Differenceing (TD)

$$V_{t+1}^\pi(s_t) = V_t^\pi(s_t) + \alpha_t \left(\underbrace{r(s_t, a_t) + \gamma V_t^\pi(s_{t+1})}_{1/(k+1)} - V_t^\pi(s_t) \right)$$

$y^{(t+1)}$
target

TD error

$$\alpha_t > 0, \quad \lim_{t \rightarrow \infty} \alpha_t = 0, \quad \sum_t \alpha_t = +\infty, \quad \sum_t \alpha_t^2 < +\infty$$

$$Q^\pi(s, a) = \mathbb{E}^{\mathbb{P}_\pi} [r(s, a) + \gamma Q^\pi(s', a') \mid s=s, A=a]$$

(s, a, r, s', a') SARSA

$$Q_{t+1}^\pi(s_{t+1}, a_{t+1}) = Q_t^\pi(s_t, a_t) + \alpha_t \left(\underbrace{r(s_t, a_t) + \gamma Q_t^\pi(s_{t+1}, a_{t+1}) - Q_t^\pi(s_t, a_t)}_{\text{TD error}} \right)$$

target

$$Q_{t+1}^{\pi}(s_t, a_t) = Q_t^{\pi}(s_t, a_t) + \alpha_t [r(s_t, a_t) + \gamma \cdot \underset{a'}{\operatorname{argmax}} Q_t^{\pi}(s_{t+1}, a') - Q_t^{\pi}(s_t, a_t)]$$

TD error

$$Q_t^{\pi}(s_{t+1}, \underset{a'}{\operatorname{argmax}} Q_t^{\pi}(s_{t+1}, a'))$$

$$= \underset{a'}{\max} Q_t^{\pi}(s_{t+1}, a')$$

$\pi(s_{t+1}) \neq \underset{a'}{\operatorname{argmax}} Q_t^{\pi}(s_{t+1}, a')$

Q-learning

$$Q(s, a) = \beta' a g(s)$$

$$\rightarrow F_{\theta}(s, a)$$

