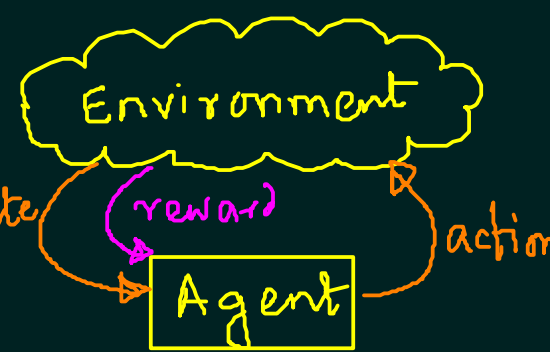
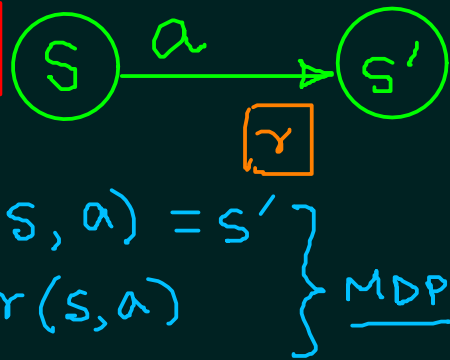


Reinforcement Learning

SUMMARY



Task: Learn an optimal policy $\pi: S \rightarrow A$

$$\left. \begin{aligned} \delta(s, a) = s' \\ r(s, a) \end{aligned} \right\} \text{MDP}$$

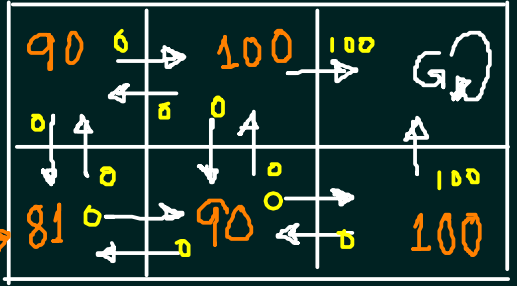
$$\pi: S_1 \xrightarrow{a_1} S_2 \xrightarrow{a_2} S_3 \xrightarrow{a_3} S_4 \rightarrow \dots$$

$$V^\pi(s_1) = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots = r_1 + \gamma V^\pi(\delta(s_1, a_1))$$

Value function

$$\pi^*(s) = \underset{a}{\operatorname{argmax}} \left[r(s, a) + \gamma V^*(\delta(s, a)) \right]$$

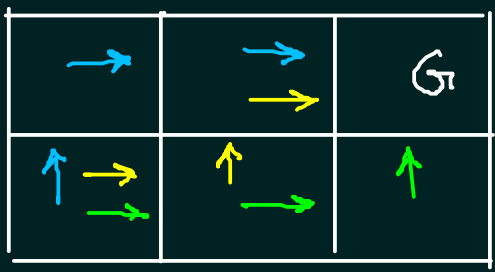
Bellman Equation



Value function compute

$$V^*(s) = 0 + 0.9 \times 90 = 81$$

Infeasible to compute in BFS manner



Q-Learning

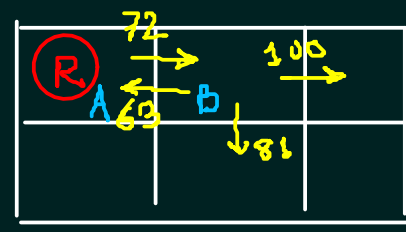
$$\pi^*(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$$

$$V^*(s) = \max_a Q(s, a)$$

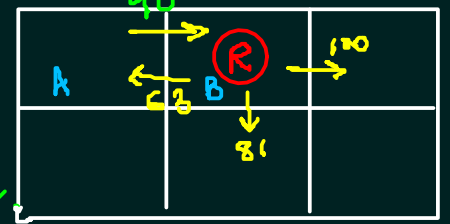
$$Q(s, a) \leftarrow r(s, a) + \gamma V^*(\delta(s, a))$$

$$= r(s, a) + \gamma \max_{a'} Q(\delta(s, a), a')$$

$$\Rightarrow \hat{Q}(s, a) = r(s, a) + \gamma \hat{Q}(s', a'), \quad s' = \delta(s, a)$$



$$\hat{Q}(A, \rightarrow) = 0 + 0.9 \max\{63, 81, 100\}$$



Exploration vs. Exploitation:

↳ Get stuck local minima ^{max.}

Solution: $P(a|s) = \frac{1}{Z} e^{\hat{Q}(s,a)/T}$

↳ alternative explore of a

— Actions are not-deterministic

$P(s' | a, s) \rightarrow$ value

$V^\pi(s_t) = E \left[\sum_{i=0}^{\infty} \gamma^i r_{t+i} \right]$ ← expected sum of cumulative rewards

$$Q(s,a) = E \left[r(s,a) + \gamma V^*(s') \right]$$

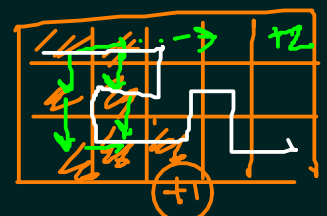
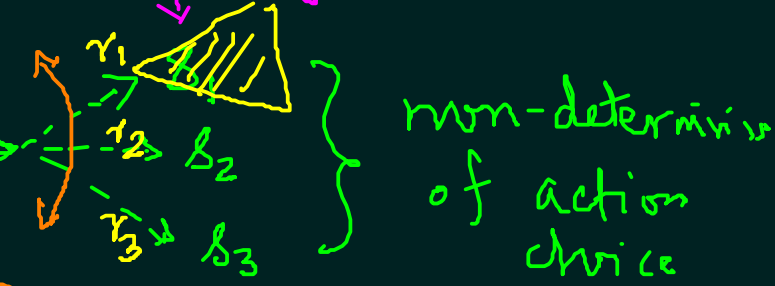
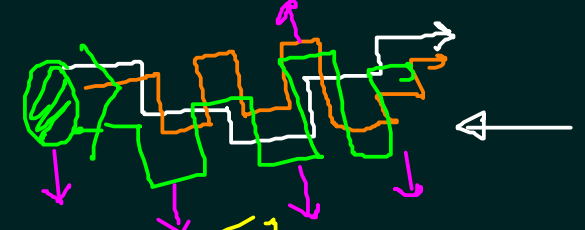
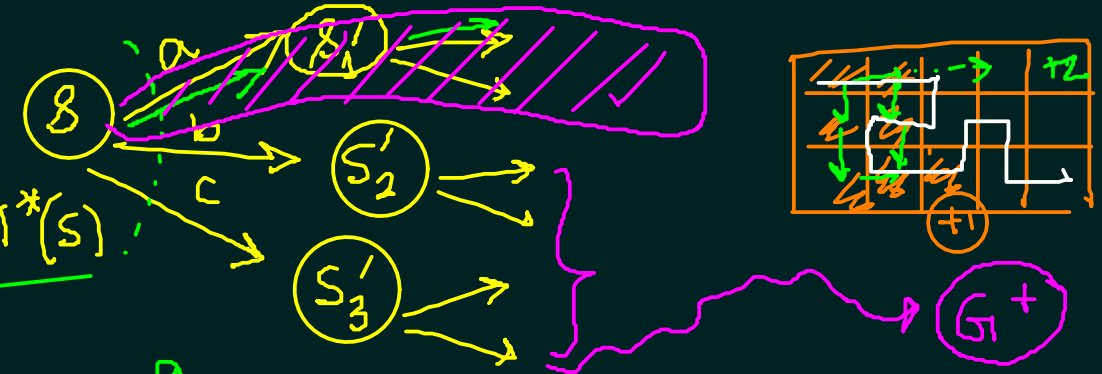
$$= E \left[r(s,a) \right] + \gamma E \left[V^*(s') \right]$$

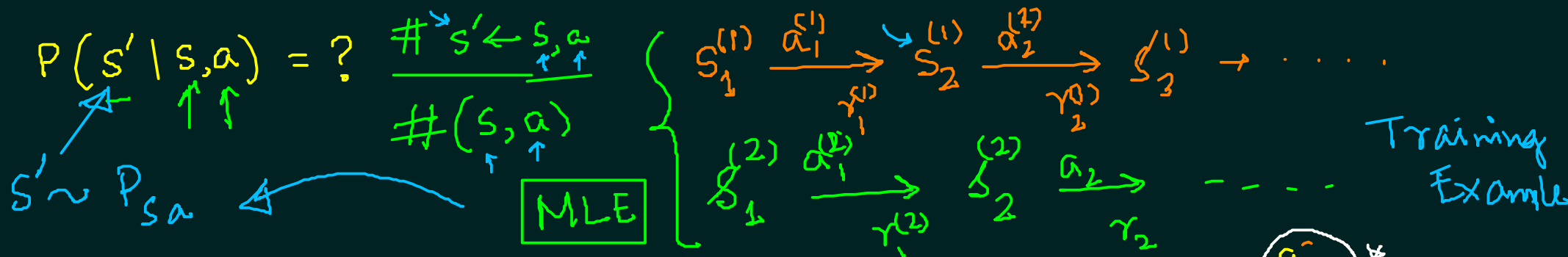
$$= E \left[r(s,a) \right] + \gamma \sum_{s'} P(s'|s,a) \cdot V^*(s')$$

$P(s'|s,a)$

— expected Q function update rule

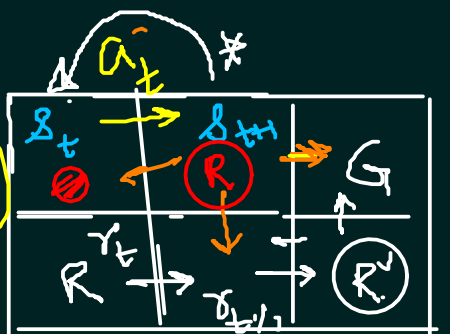
$\hat{Q}(s,a) \leftarrow E[r(s,a)] + \gamma \sum_{s'} P(s'|s,a) \max_{a'} \hat{Q}(s',a')$





Ramifications:

$\hat{Q}^{(1)}(s_t, a_t) \leftarrow r_t + \gamma \max_a \hat{Q}(s_{t+1}, a)$



$\hat{Q}^{(2)}(s_t, a_t) \leftarrow r_t + \gamma r_{t+1} + \gamma^2 \max_a \hat{Q}(s_{t+2}, a)$

← 2-step

$\hat{Q}^{(n)}(s_t, a_t) \leftarrow r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n \max_a \hat{Q}(s_{t+n}, a)$

← n-steps

$0 < \lambda < 1$

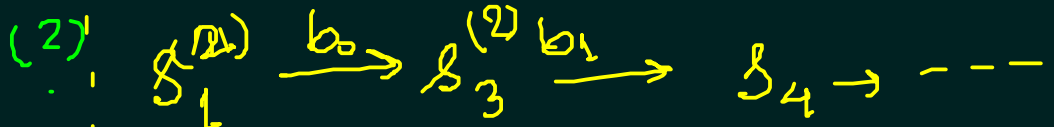
(Sutton)

$\hat{Q}^\lambda(s_t, a_t) = (1-\lambda) \left[\hat{Q}^{(1)}(s_t, a_t) + \lambda \hat{Q}^{(2)}(s_t, a_t) + \lambda^2 \hat{Q}^{(3)}(s_t, a_t) + \dots \right]$

$\Rightarrow \hat{Q}^\lambda(s_t, a_t) = r_t + \gamma \left[(1-\lambda) \max_a \hat{Q}(s_{t+1}, a) + \lambda \hat{Q}^\lambda(s_{t+1}, a_{t+1}) \right]$

Temporal Difference
 RL

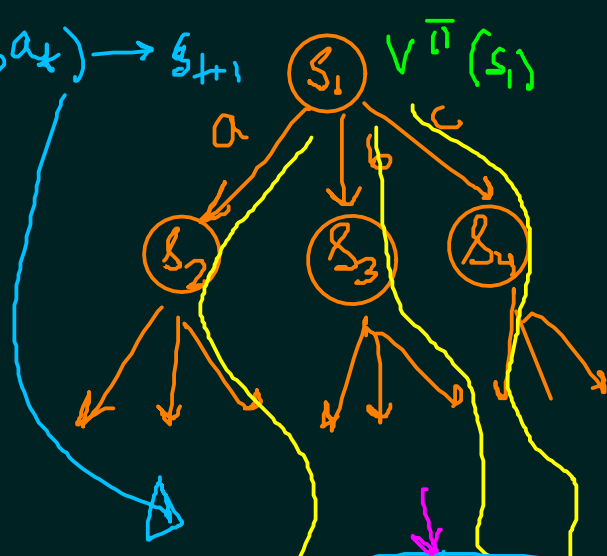
Monte-Carlo Policy Evaluation:



(k) $\pi^\pi(s) \rightarrow$ policy

$$\left. \begin{aligned} \pi(s_1) &= v_1 \\ \pi(s_2) &= v_2 \end{aligned} \right\} TE$$

$$v^\pi(s_1) = \sum \text{avg of } v_i$$



$$s_{t+1} = A s_t + B a_t$$



(x, y, z)
 (ϕ, ψ, θ) state

$(x, y, \theta) \leftarrow$ state
 $(\dot{x}, \dot{y}, \dot{\theta}) \leftarrow$ action

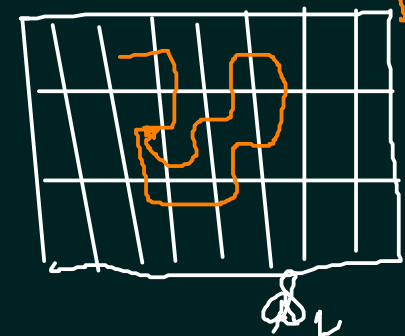
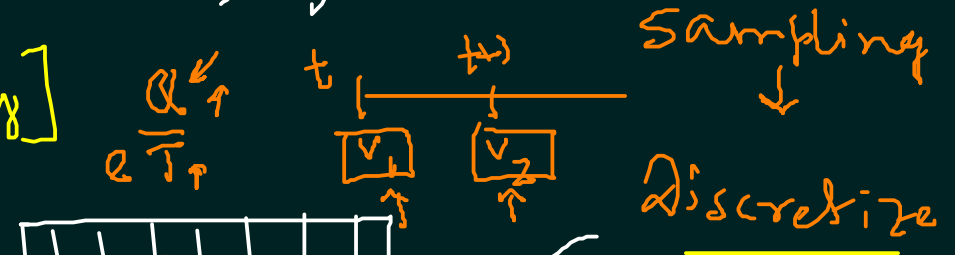
$(\ddot{x}, \ddot{y}, \ddot{z}) \leftarrow$ action
 $(\dot{\phi}, \dot{\psi}, \dot{\theta})$ action

\mathbb{R}^n [Curse of Dimensionality]

$n \rightarrow (k \text{ values}) \rightarrow k^n$



Reward Hacking
 Reward Shaping
 Safety Shield



$s_1 \rightarrow \uparrow$
 $s_2 \rightarrow \downarrow$