# Reinforcement Learning



reward $r_t$

Environment
state $S_t$

action $a_t$

Agent

$r_t = r(S_t, a_t)$

or

$r(S_1 - S_t, a_1 - a_t)$
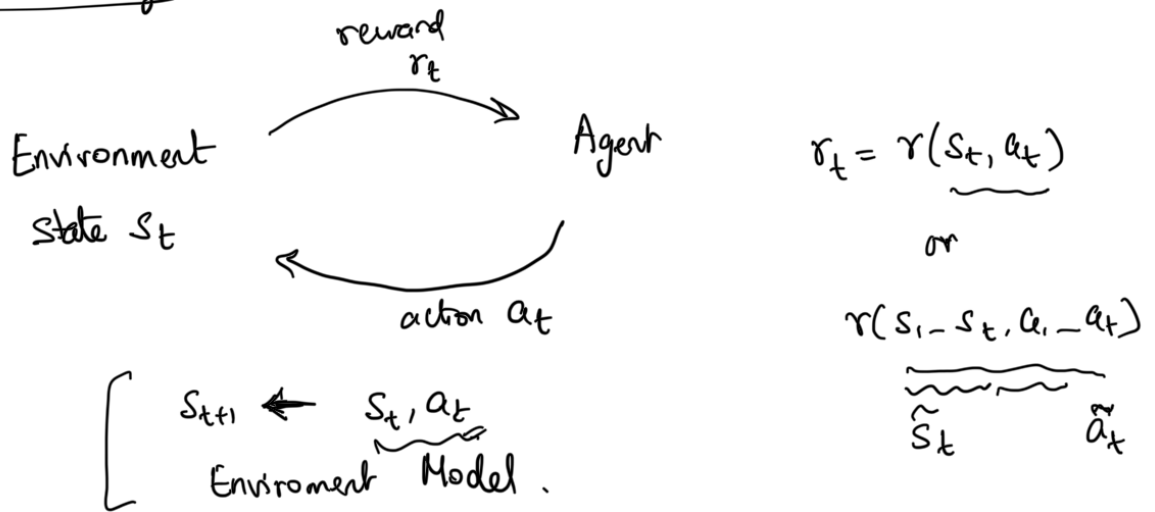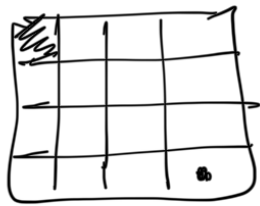$\tilde{S}_t$      $\tilde{a}_t$

$$S_{t+1} \leftarrow S_t, a_t$$
Enviroment Model.

eg. grid world



$S_t$ = location / neighborhood of agent

$a_t$ = up/down/left/right

$r_t = 1$

## Markov Decision Process

state space $S$

action space $A$

initial state distribution $d_0 \in \Delta(S)$

state transition probability $P: S \times A \rightarrow \Delta(S)$    $P(s'|s,a)$

reward function $R: S \times A \rightarrow [0, R_{max}=1] / \Delta[0,1]$

discount factor $r \in [0,1]$    how much future rewards are discounted

or

time horizon $H$

MDP $(S, A, d_0, P, R, r/H)$

Interaction protocol      Env. state $s_1 \sim d_0$

for $t = 1, 2, \dots$

agent takes action $a_t \in A$

attains rewards $r_t = R(S_t, a_t)$

observes next state of env $S_{t+1} \sim P(S_t, a_t)$

① Planning : Given state model, find desired policy

② learning : Given trajectories observation, learn desired policy.

deterministic policy $\quad \pi : S \to A \qquad a_t = \pi(S_t)$

Stochastic $\qquad \quad \sigma : S \to \Delta(A) \qquad a_t \sim \pi(S_t)$

Goal: Choose policy $\pi$ to max expected discounted sum of rewards starting at state $s_1$

$$E\left[ \sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid \pi, s_1 \right] \qquad \text{or} \qquad E\left[ \sum_{t=1}^{H} r_t \mid \pi, s_1 \right]$$

Note: $\qquad 0 < \sum_{t=1}^{\infty} \gamma^{t-1} r_t \leq R_{max} \sum_{t=1}^{\infty} \gamma^{t-1} = \frac{R_{max}}{1-\gamma} \approx R_{max} H$

$$\begin{array}{c} \text{effective} \\ \text{horizon} \end{array} \quad \frac{1}{1-\gamma} \approx H$$

Value function $\qquad V_\pi(s) = E\left[ \sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid \pi, S_1 = s \right] \qquad \forall s \in S.$

Value of policy $\pi$
starting at state $s$

Action-value or Q function

$$Q_\pi(s,a) = E\left[ \sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid \pi, S_1 = s, a_1 = a \right] \qquad \leftarrow$$

Bellman equations for policy evaluation

$$V_\pi(s) = Q_\pi(s, \pi(s))$$

$$Q_\pi(s,a) = R(s,a) + \underbrace{E\left[ \sum_{t=2}^{\infty} \gamma^{t-1} r_t \mid \pi, S_1 = s, a_1 = a \right]}$$

$$= \gamma \; E_{s' \sim P(\cdot \mid s,a)}\left[ E\left[ \sum_{t=2}^{\infty} \gamma^{t-2} r_t \mid \pi, S_2 = s' \right] \right]$$

$$= \gamma \; \underbrace{E_{s' \sim P(\cdot \mid s,a)}\left[ E\left[ \sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid \pi, S_1 = s' \right] \right]}_{V_\pi(s') \; = \; Q_\pi(s', \pi(s'))}$$

If $S, A$ finite, write as linear equations in matrix-vector form

$$\underset{|S\times A|\times 1}{Q_\pi} = \underset{|S\times A|\times 1}{R} + \gamma \underset{|S\times A|\times|S|}{\underline{P}} \underset{|S|\times 1}{V_\pi} \qquad P = P(s',a'|s)$$

Bellman eqs
$$\underset{|S\times A|\times 1}{Q_\pi} = R + \gamma \underset{|S\times A|\times|S\times A|}{P_\pi} \underset{|S\times A|\times 1}{Q_\pi} \qquad \begin{array}{l} \binom{P_\pi}{=} P(s',a'|s,a) \\ {\scriptstyle (s,a)} = P(s'|s,a) \cdot \\ {\scriptstyle (s',a')} \quad \pi(a'|s') \cdot \end{array}$$

set of $|S\times A|$ linear equations

$$Q_\pi - \gamma P_\pi Q_\pi = R$$

Bellman sol$^n$
$$Q_\pi = (I-\gamma P_\pi)^{-1} R \qquad \begin{array}{l} x^T(I-\gamma P_\pi)x \\ = x^T x - \gamma x^T P_\pi x \\ \geq (1-\gamma)\, x^T x > 0 \end{array}$$

$\rightarrow$ State-action value function is linear in $R$.

$\rightarrow$ rows of $(I-\gamma P_\pi)^{-1}$ - exp no. of times policy $\pi$ visits each state, action path

state-action occupancy distribution $\underline{d^{\pi,s}} = \underline{(I-\gamma P_\pi)^{-1}(1-\gamma)}$

$$\mathbb{1}^T(I-\gamma P_\pi)^{-1} = \mathbb{1}^T \sum_{t=1}^{\infty} \gamma^{t-1}(P_\pi)^{t-1} = \sum_{t=1}^{\infty} \gamma^{t-1} \underbrace{\mathbb{1}^T(P_\pi)^{t-1}}_{1}$$

$$= \frac{1}{1-\gamma} \mathbb{1}$$

Bellman Optimality $\qquad \pi^* = \arg\max_\pi V_\pi(s) \qquad \leftarrow$

Thm: There always exists a stationary and deterministic policy $\pi^*$ that simultaneously maximizes $V_\pi(s) \; \forall s \in S$ & $Q_\pi(s,a) \; \forall s \in S, a \in A$.

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

Bellman Optimality Eqns
$$V^*(s) = \max_{a\in A} \left[ R(s,a) + \gamma \sum_{s'\in S} P(s'|s,a) V^*(s') \right] \qquad \checkmark$$

nonlinear eqs.
$$V^* = T_v V^*$$

If we know $V^*/Q^*$ can find $\pi^*$

$$\implies \quad \pi^*(s) = \underset{a \in A}{\text{argmax}} \quad R(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) V^*(s')$$

$$Q^*(s,a) = R(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) \underset{a \in A}{\max} Q^*(s',a) \quad \checkmark$$

$$\to \quad \underline{Q^* = T_g Q^*} \quad \text{nonlinear eq's}$$

**Thm:** $Q$ is optimal iff satisfies $Q = TQ$

**Proof:** Sufficiency — by construction of $T$.

necessity — if $Q = TQ$ then $Q = Q^*$

$$\to \quad \underline{Q = TQ} = R + \gamma P_\pi Q \quad \text{where } \pi = \underset{a}{\text{argmax}} \, Q(s,a)$$
$$= (I - \gamma P_\pi)^{-1} R \qquad =: \pi_Q \quad \forall s$$

$$\to \quad [P_\pi Q - P_{\pi'} Q]_{s,a} = E_{s' \sim P(\cdot|s,a)} [Q(s', \pi(s')) - Q(s', \pi'(s'))] \geq 0$$
$$\because \pi = \pi_Q.$$

$$Q - Q_{\pi'} = (I - \gamma P_\pi)^{-1} R - (I - \gamma P_{\pi'})^{-1} R$$
$$= (I - \gamma P_{\pi'})^{-1} [(I - \gamma P_{\pi'}) - (I - \gamma P_\pi)] \underbrace{Q}_{(I - \gamma P_\pi)^{-1} R}$$
$$= \gamma (I - \gamma P_{\pi'})^{-1} \underbrace{(P_\pi - P_{\pi'})}_{\geq 0} Q$$
$$\underset{\geq 0}{\downarrow} \qquad \underset{\geq 0}{} \qquad \underset{\geq 0}{}$$

$$Q \geq Q_{\pi'} \quad \forall \pi' \qquad \qquad \oplus$$

**Non-stationary** (Time dependent) MDP . finite horizon.

$$\text{MDP} ( S, A, \{P_h\}_{h \in 1..H}, \{R_h\}_{h \in 1..H}, H)$$

$$P_h : S \times A \to \Delta(S) \qquad P_h(s'|s,a) \quad \text{at each time step } h.$$
$$R_h \qquad V_h^\pi(s) \qquad Q_h^\pi(s,a) \qquad \pi_h^* = \underset{a}{\text{argmax}} \, Q_h^{*\pi}(s,a)$$