# RL

MDP $(S, A, P, R, H/\gamma)$  $\overset{d_0}{}$

$\pi^* = \arg\max_{\pi} E[V(\pi)]$

$V(\pi) = \sum_{t=1}^{T} \sum_{h=1}^{H} r_h^t$

**Planning**: Given $\underline{P, R}, d_0$

find $\pi^*$

$d_0 : S \to R$

$R : S \times A \to R$

$P : S \times A \to S$ (drives complexity of learning)

Dynamic Prog

→ Policy iteration

→ Value iteration

$\pi$ random

$Q_\pi$ evaluate policy

$\pi_{Q_\pi}$ improvement policy

$Q_0$ random

$\pi_k = T Q_{k-1}$ improve policy

$\pi_k \leftarrow$ greedy $Q_{\pi_k}$

exp conv.     $\|Q_{\pi_k} - Q^*\| \leq \gamma^k \|Q_0 - Q^0\|$

value     $\|V_k - V^+\| \simeq \|Q_k - Q^*\| \leq \gamma^k \|Q_0 - Q^*\|$
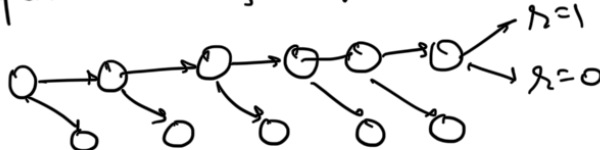
$Q_{\pi_k} \neq Q_k$

✓ Linear Prog.

# Learning Policy

Don't know state transitions $\underline{P}$ (rewards R, initial state dis" $d_0$)

Need to "explore" to learn $P, R, d_0$.

Why not use supervised learning (explore randomly?



$h=1$

$r=0$

$2^H$ episodes

Episodic MDP     $M = (S, A, P, R, H)$

starting state deterministic $S_1$

↑ reward     known

Episode $t$ generates trajectory $\tau_t = \{s_1^t, a_1^t, r_1^t \ldots s_H^t, a_H^t, r_H^t\}$ ←

$$\text{Regret }(T) = T\,E[V_{\pi^*}(s)] - E\left[\sum_{t=1}^{T}\sum_{h=1}^{H} r_h^t\right]$$

$$= \sum_{t=1}^{T}\left(E[V_{\pi^*}(s)] - E[V_{\bar\pi_t}(s)]\right)$$

↳ policy $\bar\pi_t$ deploy at time $t$.

① Natural idea

    1) estimate $\hat{P}_t$ using trajectories in episode $t$

    2) plan $\hat\pi$ using $\hat{P}$

    3) deploy $\hat\pi$ next episode

greedy approach $\equiv$ $\varepsilon$-greedy

MAB    $\mu(a^*) - \mu(a_t) \leq \underbrace{\mu(a^*) - \hat\mu_t(a^*)}_{} + \underbrace{\hat\mu_t(a^*) - \hat\mu_t(a_t)}_{\leq 0} + \underbrace{\hat\mu_t(a_t) - \mu(a_t)}_{}$

                                       greedy

                                       $a_t = \arg\max_a \hat\mu_t(a)$

$$\leq \sigma_t(a^*) - \sigma_t(a_t) \quad \text{‖}$$

            ↑         ↳ shrinks

          may not
           shrink

② UCB-VI  (Upper Confidence Bound - Value Iteration)

MAB    $\mu(a^*) - \mu(a_t) \leq \hat\mu_t(a^*) + \sigma_t(a^*) - \hat\mu_t(a_t) + \sigma_t(a_t)$

              $\underset{\text{UCB}}{\leq}$   $\hat\mu_t(a_t) + \sigma_t(a_t) - \hat\mu_t(a_t) + \sigma_t(a_t)$

                   $= \underset{\text{shrinks}}{2\sigma_t(a_t)}$   ‖       $\sigma_t(a) \sim \sqrt{\dfrac{\log T}{n_t(a)}}$

Goal:   $E[V_{\pi^*}(s)] - E[V_\pi(s)] \leq$ confidence (shrink)

    →    $\underbrace{\hspace{4cm}}$        $=$ compute using data

        Optimistic Regret Decomposition

$$E[V_{\pi^*}(s)] - E[V_{\hat\pi}(s)] \leq \sum_{h=1}^{H} E_{(s_h, a_h) \sim d_h^{\bar\pi}}\left[\text{conf}_h(s_h, a_h)\right] ←$$

  → if $\bar\pi$ is greedy policy corresponding to $\bar{Q}$ :   $\bar\pi = \arg\max_a \bar{Q}(s,a)$

    where     $Q_h^*(s,a) \leq \bar{Q}_h(s,a) \leq T\bar{Q}_{h+1}(s,a) + \text{conf}_h(s,a)$

$\forall h, s, a$      Optimistic      nearly Bellman consistent

Proof: Assuming $\bar{Q}$ exists with these properties,

$$E[V_{\pi^*}(s)] - E[V_{\bar{\pi}}(s)] = V_1^*(s_1) - V_1^{\bar{\pi}}(s_1) \qquad s_1 \text{ deterministic}$$

$$= Q_1^*(s_1, \pi^*(s_1)) - Q_1^{\bar{\pi}}(s_1, \bar{\pi}(s_1))$$

$$\underset{\text{optimistic}}{\leq} \bar{Q}_1(s_1, \pi^*(s_1)) - Q_1^{\bar{\pi}}(s_1, \bar{\pi}(s_1))$$

$$\underset{\substack{\text{ bin } \bar{\pi} \text{ is greedy} \\ \text{policy for } \bar{Q}}}{\leq} \bar{Q}_1(s_1, \bar{\pi}(s_1)) - Q_1^{\bar{\pi}}(s_1, \bar{\pi}(s_1)) \quad\text{——(a)}$$

$$\underset{\substack{\text{nearly Bellman} \\ \text{consistency of } \bar{Q}}}{\leq} T\bar{Q}_2(s_1, \bar{\pi}(s_1)) + \text{conf}_1(s_1, \bar{\pi}(s_1))$$
$$\qquad\qquad\qquad - Q_1^{\bar{\pi}}(s_1, \bar{\pi}(s_2))$$

$$=$$

$$= \underline{R(s_1, \bar{\pi}(s_1))} + E_{s_2 \sim P(\cdot|s_1, \bar{\pi}(s_1))}\left[\bar{Q}_2(s_2, \bar{\pi}(s_2))\right] + \underline{\text{conf}_1(s_1, \bar{\pi}(s_1))}$$
$$\text{Bellman operator when } \bar{\pi} \text{ greedy for } \bar{Q}$$

$$\qquad - \underline{R(s_1, \bar{\pi}(s_1))} - E_{s_2 \sim P(\cdot|s_1, \bar{\pi}(s_1))}\left[Q_2^{\bar{\pi}}(s_2, \bar{\pi}(s_2))\right]$$
$$\text{linear Bellman eq}$$

$$(b) = E_{s_2 \sim P(\cdot|s_1, \bar{\pi}(s_1))}\left[\bar{Q}_2(s_2, \bar{\pi}(s_2)) - Q_2^{\bar{\pi}}(s_2, \bar{\pi}(s_2))\right] + E_{(s_2) \sim d_1^{\bar{\pi}}}\text{conf}_1(s_2, a)$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \uparrow \text{ deterministic}$$

Notice that (a) - (b) form a recursion

$$E[V_{\pi^*}(s)] - E[V_{\bar{\pi}}(s)] \leq \sum_{k=1}^{K} E_{(s_k, a_k) \sim d_k^{\bar{\pi}}}\left[\text{conf}_k(s_k, a_k)\right]$$

$$\therefore \bar{Q}_{H+1} = 0 \neq Q_{H+1}^{\bar{\pi}} = 0 . \quad \boxed{\Box}$$

How to construct $\bar{Q}$ ?

Optimistic planning via bonus

$N_{t-1}(s, a, s')$   — # times see $(s, a, s')$ in past $t-1$ trajectories

$N_{t-1}(s, a)$   — # "    $(s, a)$

$$P^{t-1}(s'|s, a) = \frac{N_{t-1}(s, a, s')}{N_{t-1}(s, a)}$$

$$Q_H^{t-1}(s,a) = K(s,a)$$

optimistic $\rightarrow$ VI update
$$Q_h^{t-1}(s,a) = \boxed{R(s,a) + \sum_{s'} P^{t-1}(s'|s,a) \max_{a'} Q_{h+1}^{t-1}(s',a') + b^{t-1}(s,a)} \quad \text{Bellman} \atop \text{eqn.}$$

$$\underbrace{}_{\text{reward bonus}} \qquad \underline{\min(H, \boxed{\ })}$$

$\longrightarrow$ Reward bonus, $b^{t-1}(s,a) = H\sqrt{\dfrac{S\Delta}{N^{t-1}(s,a)}}$    $\Delta = \log\dfrac{SAHT}{\delta}$

$\equiv$ VI with MDP $(S, A, P^{t-1}, R+b^{t-1}, d_0, H)$

**Lemma:** $w.p \geq 1 - \delta$
$$\forall t, h, s, a \qquad Q_h^*(s,a) \leq Q_h^{t-1}(s,a) \leq TQ_{h+1}^{t-1}(s,a) + \underbrace{2b^{t-1}(s,a)}_{\text{conf.}}$$

**Proof:** By martingale version of Bernstein inequality
$$\forall s,a,t \qquad |P^{t-1}(s'|s,a) - P(s'|s,a)| \leq \sqrt{\dfrac{P(s'|s,a)\log\dfrac{SAHT}{\delta}}{N^{t-1}(s,a)}}\checkmark$$

$$+ O\left(\dfrac{1}{N^{t-1}(s,a)}\right)$$

$$\Rightarrow \quad \|P^{t-1}(s,a) - P(s,a)\|_{TV} = \underline{\sum_{s'} \tfrac{1}{2}|P^{t-1}(s'|s,a) - P(s'|s,a)|}$$

$$\leq \sqrt{\sum_{s'} 1^2 \cdot \sum_{s'} |P^{t-1}(s'|s,a) - P(s'|s,a)|^2}\checkmark$$

$$\leq \sqrt{\dfrac{S\Delta}{N^{t-1}(s,a)}\sum_{s'} P(s'|s,a)} =: \dfrac{b^{t-1}(s,a)}{H}$$

Induction to show optimistic $Q_h^{t-1}$
$$Q_{h+1}^{t-1}(s,a) \geq \underline{Q_{h+1}^*(s,a)} \quad \forall s,a$$

$$Q_h^{t-1}(s,a) = \underline{R(s,a)} + \underline{b^{t-1}(s,a)} + \sum_{s'} P^{t-1}(s'|s,a)\max_{a'} \underline{Q_{h+1}^{t-1}(s',a')}$$

$$\geq \quad " \quad + \quad " \quad + \quad \sum_{s'} P^{t-1}(s'|s,a) V_{h+1}^*(s') \qquad \text{optimism} \atop \text{of } h+1$$

$$= \quad " \quad + \quad " \quad + \quad \sum_{s'}(P^{t-1}(s'|s,a) - P(s'|s,a)) \underbrace{V_{h+1}^*(s')}_{}$$

$$+ \sum_{s'} P(s'|s,a) V_{h+1}^b(s')$$

$$\geq \quad " \quad + \quad " \quad - \|P^{t-1}(s,a) - P(s,a)\|_{TV} \cdot H \quad + \quad "$$

$$\geq \quad R(s,a) + \sum_{s'} P(s'|s,a) V_{h+1}^*(s') = Q_h^*(s,a)$$

$$Q_h^{t+1}(s,a) = R(s,a) + b^{t+1}(s,a) + \sum_{s'} P^-(s'|s,a) \max_{a'} Q_{h+1}^{t+1}(s',a')$$

$$\leq b^{t+1}(s,a) + T Q_{h+1}^{t+1}(s,a) + \underbrace{\sum_{s'} \underbrace{(P^{t-1}(s'|s,a) - P(s'|s,a))}_{\cdot \max_{a'} Q_{h+1}^{t+1}(s',a') \frac{b^{t+1}}{H}}}_{H}$$

$$\leq \underbrace{2 b^{t+1}(s,a)}_{conf} + T Q_{h+1}^{t+1}(s,a)$$

## Regret of UCB-VI

### Optimistic Regret Decomp.

$$\bar{Q}_h \equiv Q_h^{t+1} \qquad \text{optimistic \& nearly Bellman consistent}$$

$$Regret(T) = \sum_{t=1}^{I} \mathbb{E}[V_{\pi^*}(s) - V_{\pi_t}(s)] \leq \sum_{t=1}^{I} \sum_{h=1}^{H} 2 \, \mathbb{E}_{s,a \sim d_h^{\pi_t}} [b^{t+1}(s,a)]$$

$$\downarrow$$
$$\text{greedy for } Q_h^{t+1}$$

$$\sum_{t=1}^{I} \sum_{h=1}^{H} H \sqrt{\frac{S\Delta}{N^{t+1}(s_h, a_h)}} = H\sqrt{S\Delta} \sum_{t,h} \frac{1}{\sqrt{N^{t+1}(s_h, a_h)}}$$

$$= H\sqrt{S\Delta} \sum_{h} \sum_{s,a} \sum_{i=1}^{N_h^T(s,a)} \frac{1}{\sqrt{i}} \qquad\qquad \sum_{i=1}^{h} \frac{1}{\sqrt{i}} \leq 2\sqrt{h}$$

$$\neq H\sqrt{S\Delta} \sum_{h} \sum_{s,a} 1 \cdot \sqrt{N_h^T(s,a)}$$

$$\leq H\sqrt{S\Delta} \sum_{h} \sqrt{\underbrace{\sum_{s,a} 1^2}_{SA} \underbrace{\sum_{s,a} N_h^T(s,a)}_{T}}$$

$$\varepsilon = H^2 S \sqrt{AT\Delta} \qquad\qquad \boxed{2}$$

$$T \cong S^2 A \qquad\qquad\qquad P: S \times A \to S$$

Can be improved:
$$\Rightarrow \quad T \cong SA \qquad\qquad \text{Optimizing Regret} \quad\Big\} \begin{array}{l}\text{easier than} \\ \text{learn} \\ \text{model.}\end{array}$$
$$\text{Find optimal policy}$$