# Local planning / Online planning

## Why / What



$A \leftarrow$ random

query

$\boxed{s, a}$

$S', r_a(s)$

$Sim$ (W)

$P$

$s_0$

$S' \sim P_a(s)$ } stochasticity leveraged to save on compute?!
$\uparrow$ random

---

$$\pi(a|s_0) = \mathbb{P}_{s_0}(A=a)$$ $\boxed{\text{Certainty Equivalence}}$

No caching : memoryless planner $\Rightarrow$ memoryless policy

$\boxed{Sim \overset{?}{=} W}$
sensitivity

Goal: ① $\underbrace{v^\pi}_{v^\pi(s)} \geq v^* - \delta 1$

$\delta > 0$: input to the planner

② efficiency

\* computation-time

\*\* #query / query-cost

$$v_0 = 0$$

$$v_{k+1} = T v_k, \quad k = 0, 1, \dots$$

$$v_k = T^k 0 \to v^*$$

$\delta$ - support policy?

$k = ?$

$k = H_{0, (1-\gamma)^2 \delta}$

$\approx \dfrac{\log(1/\delta)}{1-\gamma}$

Sampling

$$\arg\max_a \ r_a(s_0) + \gamma \langle P_a(s_0), v^* \rangle$$

optimal actions!

$$S_1, \dots, S_m \sim P_a(s_0)$$

$$\langle \hat{P}_a(s_0), v^* \rangle = \frac{1}{m} \sum_{i=1}^m v^*(S_i)$$

$$\arg\max_a \ r_a(s_0) + \gamma \langle \hat{P}_a(s_0), v_k \rangle$$

Independent of the size of the state space !

$\leftarrow$ $O(A^k)$

Deterministic MDP          $s' = g(s, a)$

$$v_k(s) = (T^k 0)(s)$$

$$= \max_a \ r_a(s) + \gamma \langle P_a(s), T^{k-1} 0 \rangle$$

$\delta_{g(s,a)}$

$$= \max_a \left( r_a(s) + \gamma \, (T^{k-1} 0)( g(s,a)) \right)$$

$\underbrace{\qquad}_{v_{k-1}}$

def $v(k, s) \# v_k(s)$
if $k = 0$ return $0$;
$q = [ r_a(s) + \gamma\, v(k-1, g(s,a)) \ \text{for } a \in A]$
return $\max(q)$

queries

$A = 3$

$\searrow O((mA)^k)$

How big should be m?

$$\text{argmax}_a \underbrace{\overline{r_a(s) + \gamma \langle P_a(s), v^* \rangle}}$$

$$\underbrace{q^*(s,a)}_{\text{"optimal value of a"}}$$

$$v^*(s) = \max_a q^*(s,a) \quad \forall s$$
$$[B.O.E.]$$

$$q^\gamma(s,a) = r_a(s) + \gamma \langle P_a(s), \max_{a'} q^*(\cdot, a') \rangle$$

$$M: \mathbb{R}^{SA} \to \mathbb{R}^S$$
$$q \mapsto (Mq)(s) = \max_a q(s,a)$$

$$q^*(s,a) = r_a(s) + \gamma \langle P_a(s), Mq^* \rangle \quad \forall_{s,a}$$
$$\Big\lfloor_{\max:}$$

B.O.E $q^*$

$$q^* = r + \gamma P M q^*$$

$$P: \mathbb{R}^S \to \mathbb{R}^{SA}$$
$$v \mapsto (Pv)(s,a) = \langle P_a(s), v \rangle$$

$$r: \mathbb{R}^{SA} \to \mathbb{R}$$
$$r(s,a) = r_a(s).$$

$$\xrightarrow{\quad} M_\pi$$

$$\tilde{T} q = r + \gamma P M q$$

$$q^* = \tilde{T} q^*$$

$$T \doteq \tilde{T} \mid q^* = T q^*$$

$$\underset{a}{\text{argmax}} \ (T^k 0)(s_0, a)$$

$$(Tq)(s,a) = r_a(s) + \gamma \underbrace{\langle P_a(s), Mq \rangle}_{\text{costly!}}$$

$$P_a(s) \Rightarrow \hat{P}_a(s)$$

$$C(s,a) = \left[ S_{sa}^{(1)}, \ldots, S_{sa}^{(m)} \right]$$

$$S_{sa}^{(i)} \sim P_a(s), \quad \begin{matrix} i = 1 \ldots m \\ \text{i.i.d.} \end{matrix}$$

$$(\hat{T}q)(s,a) = r_a(s) + \gamma \frac{1}{m} \sum_{s' \in C(s,a)} (Mq)(s')$$

$$= r_a(s) + \frac{\gamma}{m} \sum_{s' \in C(s,a)} \underset{a'}{\max} \ q(s',a)$$

$$T \approx \hat{T} \qquad \text{``random approx''}$$

$$\boxed{\underset{a}{\text{argmax}} \ (\hat{T}^k 0)(s_0, a)}$$
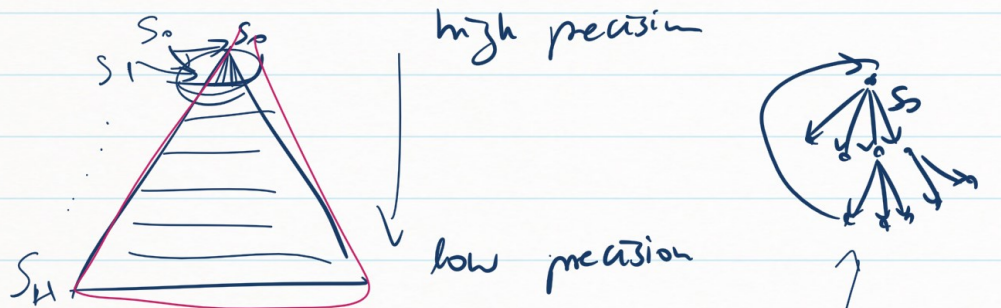
Complexity?   Branching: $mA$

Depth $k$: $O((mA)^k)$

cost / query - cost

memorization to get $C(s,a)$

$\hat{T}^H O \approx q^*$  only true at $s_0$



high precision

low precision

$S_0 = \{s_0\}$

$S_1 = \{s_0\} +$ neighbors

$\vdots$

$S_h = \{ s \in S \mid dist(s_0, s) \leq h \}$

$S_H =$ all the states encountered.

---

$|(\hat{T}^H O)(s_0, a) - q^*(s_0, a)| \leq ?$

$\rightarrow$ shell

$\delta_H = \| \hat{T}^H O - q^* \|_{S_0}$

$\delta_{H-1} = \| \hat{T}^{H-1} O - q^* \|_{S_1}$

$\boxed{\begin{array}{l} \|q\|_S = \\ = \max |q(s,a)| \\ \quad s \in S \\ \quad a \in A \end{array}}$

$\vdots$

$\delta_h = \| \hat{T}^h O - q^* \|_{S_{H-h}}$

$\vdots$

$\rightarrow \delta_0 = \| \underbrace{\hat{T}^0 O}_{0} - q^* \|_{S_H} \leq \frac{1}{1-\gamma}$

$h > 0:$

$\delta_h = \| \hat{T}^h O - q^* \|_{S_{H-h}}$

$\leq \underbrace{\| \hat{T}^h O - \hat{T} q^* \|_{S_{H-h}}}_{\leq \gamma \, \delta_{h-1}} + \boxed{\| \hat{T} q^* - T q^* \|}_{S_H}$

$+ \cdot \frac{\varepsilon!}{1-\gamma}$

$$\|\hat{T}^h 0 - \hat{T} q^* \|_{S_{H-h}}$$

$$= \| \hat{T} \underbrace{\hat{T}^{h-1} 0}_{u} - \underbrace{\hat{T} q^*}_{v} \|_{\boxed{S_{H-h}}}$$

$$u' = Mu$$
$$v' = Mv$$

$$(\hat{T} u)(s,a) = r_a(s) + \frac{\gamma}{m} \sum_{s' \subset C(s,a)} u'(s')$$

$$dist(s_0, s) \leq H-h$$

$$dist(s_0, s') \leq \underline{H-h+1}$$

$$\forall s' \in C(s,a)$$

$$\|\hat{T} u - \hat{T} v \|_{S_{H-h}} \leq$$

$$\leq \max_{\substack{s \in S_{H-h} \\ a \in A}} \left| \frac{\gamma}{m} \sum_{s' \subset C(s,a)} u'(s') - v'(s') \right|$$

$$\leq \gamma \max_{s' \in S_{H-h+1}} |u'(s') - v'(s')|$$

$$\leq \gamma \underbrace{\|u - v\|_{S_{H-h+1}}}_{\delta_{h-1}}$$

$$\delta_h \leq \gamma \delta_{h-1} + \frac{\varepsilon'}{1-\gamma} \qquad \boxed{\varepsilon'} ?$$

$$\vdots$$

$$\delta_H \leq \text{small} \cdots$$

$$(m, H) = f(\delta)$$
$$?$$

finish