

Meta-Algorithm (Approx. Dynamic Programming)

initialize π^0

for $t=0, 1, \dots$

1) $\hat{Q}^t = \text{SAMPLEANDEVAL}(\pi^t)$

2) $\pi^{t+1} = \text{IMPROVEMENT}(\hat{Q}^t)$

1) Supervision via Rollouts

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid \begin{matrix} s_0, a_0 = s, a \\ P, \pi \end{matrix} \right]$$

Alg: Infinite Rollout:

$s_0 = s, a_0 = a$

for $t=0, 1, \dots$

take action a_t , observe r_t & $s_t \sim P(s_t, a_t)$

update $a_{t+1} = \pi(s_{t+1})$

return $y = \sum_{t=0}^{\infty} \gamma^t r_t$

Then $\mathbb{E}(y) = Q^\pi(s, a)$

Alg: Rollout:

$s_0 = s, a_0 = a$

for $t=0, 1, \dots$

take action a_t , observe r_t & $s_{t+1} \sim P(s_{t+1}, a_t)$

with probability $1-\gamma$:

Break and return $y = \sum_{k=0}^t r_k$

update $a_{t+1} = \pi(s_{t+1})$

$\mathbb{E}[y]$?
Alg

$$y = \begin{cases} r_0 & \text{w.p. } 1-\gamma \\ r_0 + r_1 & \text{w.p. } \gamma(1-\gamma) \\ r_0 + r_1 + r_2 & \text{w.p. } \gamma^2(1-\gamma) \\ \vdots \\ \sum_{k=0}^t r_k & \text{w.p. } \gamma^t(1-\gamma) \end{cases}$$

$$\mathbb{E}(y) = r_0(1-\gamma) \left(\sum_{k=0}^{\infty} \gamma^k \right) + r_1(1-\gamma)\gamma \left(\sum_{k=0}^{\infty} \gamma^k \right) + \dots = \sum_{k=0}^{\infty} \gamma^k r_k$$

y is an unbiased estimate of $Q^\pi(s, a)$
 $\{(s_i, a_i, y_i)\}_{i=1}^N$

Recall: prediction error guarantee in supervised ML
 $x \sim \mathcal{D}_x \quad y = f_*(x) + w \quad (\mathbb{E}(y) = f_*(x))$

$$\mathbb{E}_{x \sim \mathcal{D}_x} [(f_*(x) - \hat{f}(x))^2] \leq \varepsilon \quad \text{(usually } \varepsilon \approx O(\frac{1}{\sqrt{N}})\text{)}$$

connection $(s, a) = \text{features } x \quad y \text{ label } (Q^\pi(s, a) = f_*)$
 use discount state-action $d_{y_0}^\pi$!

Alg: Sample

initialize $s_0 \sim \mathcal{Y}_0 \quad a_0 = \pi(s_0)$

for $t=0, 1, \dots$

Take action a_t , observe $s_{t+1} \sim P(s_t, a_t)$

w.p. $1-\delta$

break & return a_t, s_t

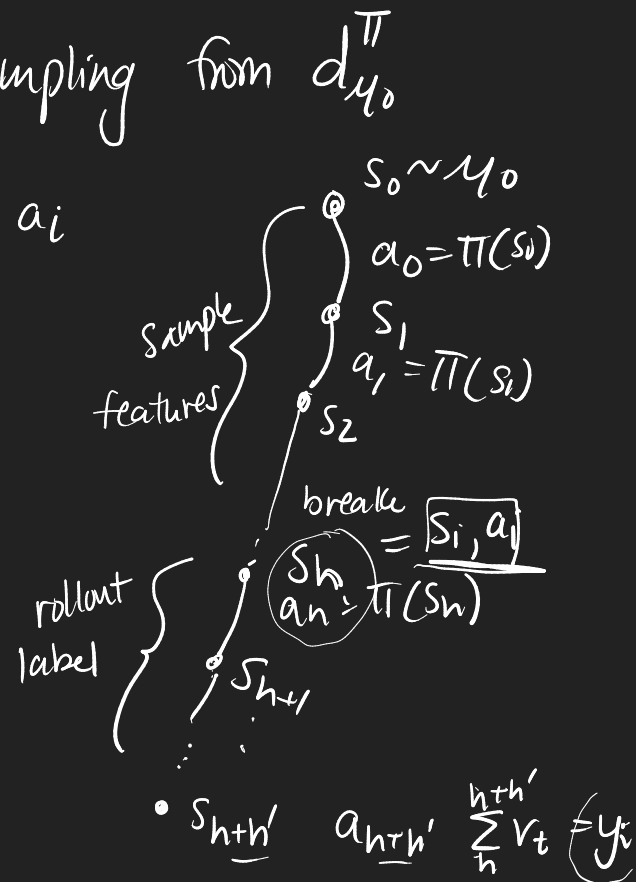
update $a_{t+1} = \pi(s_{t+1})$

This algorithm is equivalent to sampling from $d_{y_0}^\pi$

To collect dataset, 1) sample s_i, a_i
 2) rollout y_i

Terminology: MCMC
 Markov Chain Monte Carlo

$$\{(s_i, a_i, y_i)\}_{i=1}^N$$



2) Approximate Policy Iteration

Putting together sampling method & estimation:

Alg: ROLLOUT EVAL(π)

for $i=1, \dots, N$

$$s_i, a_i = \text{SAMPLE}(\pi)$$

$$y_i = \text{ROLLOUT}(s_i, a_i, \pi)$$

$$\hat{Q}^\pi = \arg \min_{Q \in \mathcal{Q}} \sum_{i=1}^N (Q(s_i, a_i) - y_i)^2$$

↑
function e.g. linear neural network $Q(s, a) = \Theta^T \phi(s, a)$

Alg: APPROX POLICY ITERATION

initialize π^0

for $t=0, 1, \dots$

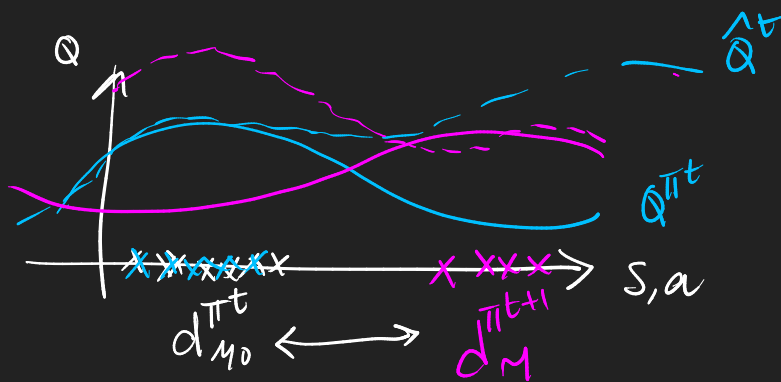
$$\hat{Q}^t = \text{ROLLOUT EVAL}(\pi^t)$$

$$\pi_{t+1}(s) = \arg \max_a \hat{Q}^t(s, a)$$

← sample & regression
← policy improvement

For policy iteration $V^{\pi^{t+1}}(s) \geq V^{\pi^t}(s)$ monotonic improvement
 $Q^{\pi^{t+1}}(s, a) \geq Q^{\pi^t}(s, a)$

$\mathbb{E}_{s, a \sim d_{\pi^t}^s} (\hat{Q}^t(s, a) - Q^{\pi^t}(s, a))^2 < \epsilon$ is our regression guarantee



oscillation!

4) Conservative Policy Iteration

Alg: conservative PI

initialize π^0

for $t=0, 1, \dots$

$$\hat{Q}^t = \text{POLUTEVAL}(\Pi^t)$$

$$\Pi'(s) = \underset{a}{\operatorname{argmax}} \hat{Q}^t(s, a)$$

$$\Pi^{t+1}(a|s) = (1-\alpha)\Pi^t(a|s) + \alpha \Pi'(a|s) \quad \forall a, s$$

$\Pi(a|s) \leftarrow$ stochastic policy
 \downarrow 1 if $a = \operatorname{argmax} \hat{Q}^t$
0 o.w.

randomize interpolation:
follow Π^t w.p. $1-\alpha$
and Π' w.p. α

α is a step size

$$\Pi^{t+1}(a|s) = \Pi^t(a|s) + \alpha (\Pi'(a|s) - \Pi^t(a|s))$$

$$d_{M_0}^{\Pi}(s, a) = \sum_{t=0}^{\infty} \gamma^t P_t^{\Pi}(s, a; M_0)$$

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1}$$

Prob. of $S_t = s$ at $t=a$ under Π $S_0 \sim M_0$