

Policy iteration:

- Start with π_0 arbitrary;

for $k=0, \dots$

• Compute Q^{π_k} (Can be computed analytically)

• Update $\pi_{k+1} = \pi_{Q^{\pi_k}}$ $(1-\gamma)(I-\gamma P^{\pi_k})^{-1} r$

Recall $P^{\pi}(s,a)(s',a') = \mathbb{P}[s'|s,a] \pi(a'|s')$.

∴ Policy evaluation followed by policy improvement.

Lemma 1 1) $Q^{\pi_{k+1}} \geq TQ^{\pi_k} \geq Q^{\pi_k}$

2) $\|Q^{\pi_{k+1}} - Q^*\|_{\infty} \leq \|Q^{\pi_k} - Q^*\|_{\infty}$

Observe:

our policies are deterministic.

$$\therefore V^{\pi_k}(s) = Q^{\pi_k}(s, \pi_k(s)) \quad \forall k, \forall s.$$

$$\therefore TQ^{\pi_k}(s, a) = (1-\gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a'} Q^{\pi_k}(s', a') \right]$$

$$\geq (1-\gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[Q^{\pi_k}(s', \pi_k(s')) \right]$$

$$TQ^{\pi_k}(s, a) = Q^{\pi_k}(s, a).$$

Now:

$$Q^{\pi_{k+1}}(s, a) = (1-\gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[Q^{\pi_{k+1}}(s', \pi_{k+1}(s')) \right]$$

$$\geq (1-\gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[Q^{\pi_k}(s', \pi_{k+1}(s')) \right]$$

$$= (1-\gamma)r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a'} Q^{\pi_k}(s', a') \right]$$

$$= TQ^{\pi_k}(s, a).$$

Now:

$$\begin{aligned}\|Q^* - Q^{\pi_k}\|_\infty &\leq \|Q^* - TQ^{\pi_k}\|_\infty = \|TQ^* - TQ^{\pi_k}\|_\infty \\ &\leq \gamma \|Q^* - Q^{\pi_k}\|_\infty \\ &\leq \gamma^k \|Q^* - Q_0\|_\infty \leq e^{-(1-\gamma)k}\end{aligned}$$

$$\begin{aligned}\therefore \forall k \geq \frac{\log(1/\epsilon)}{1-\gamma} &\leq \exp^{-\log(1/\epsilon)} \\ &\leq \epsilon \\ &= \epsilon\end{aligned}$$

Thm:

For any two policies π, π' , we have:

$$\left(\forall s \in S, \mathbb{E}_{a \sim \pi'(s)} [Q_\pi(s, a)] \geq \mathbb{E}_{a \sim \pi(s)} [Q_\pi(s, a)] \right) \Rightarrow$$

$$\left(\forall s, V_{\pi'}(s) \geq V_\pi(s) \right)$$

Note - $Q_\pi(s, a)$

Proof:

$$V_{\pi}(s) - (1-\gamma) \mathbb{E}_{a \sim \pi(s)} Q_{\pi}(s, a)$$

$$\leq (1-\gamma) \mathbb{E}_{a \sim \pi'(s)} Q_{\pi}(s, a)$$

$$= (1-\gamma) \mathbb{E}_{a \sim \pi'(s)} \left[r(s, a) + \gamma V_{\pi}(s_1) \mid s_0 = s \right]$$

$$= (1-\gamma) \mathbb{E}_{a \sim \pi'(s)} \left[r(s, a) + \gamma \mathbb{E}_{a_1 \sim \pi(s_1)} (Q_{\pi}(s_1, a_1)) \mid s_0 = s \right]$$

$$\leq (1-\gamma) \mathbb{E}_{a \sim \pi'(s)} \left[r(s, a) + \gamma \mathbb{E}_{a_1 \sim \pi'(s_1)} (Q_{\pi}(s_1, a_1)) \mid s_0 = s \right]$$

$$= (1-\gamma) \mathbb{E}_{a \sim \pi'(s)} \left[r(s, a) + \gamma r(s_1, a_1) + \gamma^2 V_{\pi}(s_2) \mid s_0 = s \right]$$

$$\Rightarrow V_{\pi}(s) \leq (1-\gamma) \mathbb{E}_{a \sim \pi'(s_t)} \left[\sum_{t=0}^T \gamma^t \mathbb{E} (r(s_t, a_t)) + \gamma^{T+1} V_{\pi}(s_{T+1}) \mid s_0 = s \right]$$

But $V_{\pi}(s_{T+1})$ is bounded $\therefore \gamma^{T+1} V_{\pi}(s_{T+1}) \rightarrow 0$ as $T \rightarrow \infty$

Note: We can ignore $(1-\gamma)$ here, or ignore it here.

\therefore Taking limit $\Rightarrow T \rightarrow \infty$,

$$V_{\pi}(s) = (1-\gamma) \mathbb{E}_{a_t \sim \pi(\cdot|s_t)} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{E} [r(s_t, a_t) | s_0 = s] \right]$$

$$= V_{\pi'}(s);$$

We get.

THM: (Bellman's optimality condition)

A policy π is optimal iff for any pair $(s, a) \in S \times A$ with $\pi(s)(a) > 0$

(recall $\pi(s)$ is a distribution on actions)

the following holds

$$a \in \operatorname{argmax}_{a' \in A} Q_{\pi}(s, a')$$

Proof: Clear.

If this condition does not hold for some (s, a) then π is not optimal.

Consider π' s.t. $\pi'(s') = \pi(s)$ for $s' \neq s$ and $\pi'(s)$ is concentrated on $\operatorname{argmax}_{a' \in A} Q_{\pi}(s, a')$.

For π' , $V_{\pi'}(s) \geq V_{\pi}(s) \forall s$.

← Conversely

If π' is not optimal, $\exists \pi$ with $V_{\pi}(s) > V_{\pi'}(s)$ some s .

But then $\exists s$ with

$$E_{a \sim \pi'(s)} [Q_{\pi}(s, a)] < E_{a \sim \pi(s)} [Q_{\pi}(s, a)].$$

Thm: \exists an optimal det policy.

SAMPLE COMPLEXITY with a GENERATIVE MODEL.

• What is the sample complexity of estimating Q^* ?

- Will assume that the reward function is known & deterministic.

- Assume we have access to a generative model

- Given s, a provides a sample $s' \sim P(\cdot | s, a)$.

- Invoke our simulator N times \forall $S \times A$ pair.

Defint:

$$\hat{P}(s' | s, a) = \frac{\#(s', s, a)}{N}$$

times our simulator
transits to s' from s on a .

Let \hat{M} = empirical MDP.

↓
uses \hat{P} instead of P .

• First Hoeffding's inequality:

• Recall Markov's inequality:

For any nonnegative random variable X

$$\Pr[X \geq t] \leq \frac{E(X)}{t}.$$

In fact: strictly
if ϕ is a monotonically \uparrow non-negative-valued
function then for any random variable X
and $t \in \mathbb{R}$,

$$\mathbb{P}_X[X > t] = \mathbb{P}_X[\phi(X) > \phi(t)] \leq \frac{\mathbb{E}[\phi(X)]}{\phi(t)}.$$

=
Using $\phi(x) = x^2$:

$$\mathbb{P}_X[|X - \mathbb{E}(X)| > t] = \mathbb{P}_X[(X - \mathbb{E}(X))^2 > t^2]$$

$$\leq \frac{\mathbb{E}[(X - \mathbb{E}(X))^2]}{t^2}$$

$$= \frac{\text{Var}(X)}{t^2} \quad \dots \text{Chebyshev.}$$

• Taking $\phi(x) = e^{sx}$, s arbitrary +ve number

$$\mathbb{P}_s [X > t] = \mathbb{P}_s [e^{sX} > e^{st}] \leq \frac{\mathbb{E}[e^{sX}]}{e^{st}}$$

find $s > 0$ which makes the RHS small!

CHEBNOFF

• Sums of independent random variables:

$S_n = X_1 + \dots + X_n$, X_i independent real-valued random variables.

$$\begin{aligned} \mathbb{P}_s \left[|S_n - \mathbb{E}(S_n)| > t \right] &\leq \frac{\text{Var}(S_n)}{t^2} \\ &= \frac{\sum \text{Var}(X_i)}{t^2} \end{aligned}$$

Writing $\sigma^2 = \frac{1}{n} \sum \text{Var}(X_i)$ we get

Then

$$\mathbb{P}_s \left[\left| \frac{S_n}{n} - \mathbb{E}(X_i) \right| \geq \epsilon \right] \leq \frac{\sigma^2}{n\epsilon^2}$$

• Not very satisfactory.

∴ Central limit theorem:

$$\mathbb{P}_0 \left[\sqrt{\frac{n}{\sigma^2}} \left(\frac{1}{n} \sum x_i - \mathbb{E}(x_i) \right) \geq y \right] \rightarrow$$
$$1 - \phi(y) \leq \frac{1}{\sqrt{2\pi}} \frac{e^{-y^2/2}}{y}$$

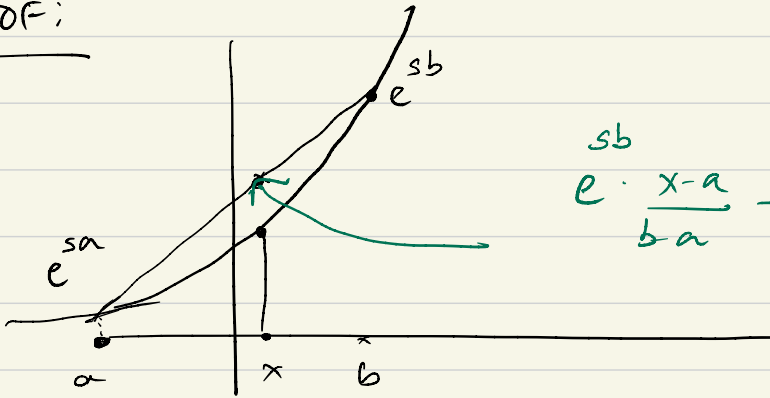
Expect: $\mathbb{P}_0 \left[\frac{1}{n} \sum x_i - \mathbb{E}(x_i) \geq \varepsilon \right] \approx$

Much smaller than
what Chebyshev
gives! $\rightarrow \exp\left(-\frac{n\varepsilon^2}{2\sigma^2}\right)$

Hoeffding: Let X be a r.v with $\mathbb{E}(X) = 0$,

$$a \leq X \leq b; \quad \forall s > 0, \quad \mathbb{E}[e^{sX}] \leq e^{s^2(b-a)^2/8}$$

PROOF:



$$e \cdot \frac{x-a}{b-a} + e \cdot \frac{b-x}{b-a}$$

$$1. \quad e^{sx} \leq \frac{x-a}{b-a} e^{sb} + \frac{b-x}{b-a} e^{sa}$$

$$\therefore \mathbb{E}[e^{sx}] \leq \frac{\mathbb{E}(x-a)}{b-a} e^{sb} + \frac{\mathbb{E}(b-x)}{b-a} e^{sa}$$

$$= \frac{b}{b-a} e^{sa} - \frac{a}{b-a} e^{sb}$$

$$\text{Setting } p = \frac{a}{b-a} \quad \therefore \quad 1-p = 1 + \frac{a}{b-a} = \frac{b}{b-a}$$

$$= (1-p) e^{-sp(b-a)} + p \cdot e^{s(b-a)(1+p)}$$
$$= \left[1-p + p e^{s(b-a)} \right] e^{-sp(b-a)}$$

Sol

$$u = s(b-a)$$

$$\mathbb{E}[e^{sx}] \leq e^{\phi(u)},$$

$$\phi(u) = -pu + \log(1-p+pe^u).$$

$$\phi'(u) = -p + \frac{pe^u}{1-p+pe^u}$$

$$= -p + \frac{p}{p+(1-p)e^{-u}}$$

$$\text{Now } \phi(0) = 0; \quad \phi'(0) = 0$$

$$\phi''(u) = \frac{p(1-p)e^{-u}}{(p+(1-p)e^{-u})^2} \leq \frac{1}{4}$$

$$\theta \in [0, \eta]$$

$$\therefore \phi(\eta) = \phi(0) + \eta \phi'(0) + \frac{\eta^2}{2} \phi''(\theta)$$

$$\leq \frac{\eta^2}{8} = \frac{s^2 (b-a)^2}{8}$$

$$\therefore \mathbb{P} \left[S_n - \mathbb{E}(S_n) \geq t \right]$$

$$\leq e^{-st} \frac{e^{\sum_{i=1}^n (x_i - \mathbb{E}(x_i))}}$$

$$= e^{-st} e^{sx_1} e^{sx_2} \dots e^{sx_n}$$

$$\leq e^{-st} \prod e^{\frac{s^2 (b_i - a_i)^2}{8}}$$

$$= e^{-st} \frac{s^2}{e} \cdot e^{\sum_{i=1}^n \frac{(b_i - a_i)^2}{8}}$$

$$\text{Choose } s = \frac{4t}{\sum_{i=1}^n (b_i - a_i)^2}$$

$$\leq e^{-\frac{4t^2}{\sum_{i=1}^n (b_i - a_i)^2}}$$

Thm: Hoeffding's tail inequality.

Let X_1, \dots, X_n be independent r.v. s.t.

$X_i \in [a_i, b_i]$ with prob. 1. Then $\forall t > 0$

$$\mathbb{P}_x [S_n - \mathbb{E}(S_n) \geq t] \leq e^{-\frac{2t^2}{\sum (b_i - a_i)^2}}$$

$$\mathbb{P}_x [S_n - \mathbb{E}(S_n) \leq -t] \leq e^{-\frac{2t^2}{\sum (b_i - a_i)^2}}$$

Coming back to the empirical model:

$$\hat{P}(s' | s, a) = \frac{m(s', s, a)}{n}$$

$n = \#$ calls to the model per (s, a) pair.

$\therefore \#$ calls $n |S||A|$.

Based on \hat{P} we get

$$\hat{T} \mathcal{J}(z) = \mathcal{J}(z) + \mathcal{J}(\hat{P} V)(z), \quad V(z) = \max_{a \in A} \mathcal{J}(z, a)$$

Can also define \hat{T}^π - operator on \mathcal{Q} as

$$\hat{T}^\pi \mathcal{Q}(z) = r(z) + \gamma \hat{P}^\pi \mathcal{Q}(z).$$

- Can define the above operators for V as well.

Algorithm: Model-based Q-value iteration.

Input: \mathcal{Q}_0 ; n - samples per $s \times A$ pair;
 k - # iterations.

\hat{P} = ESTIMATE MODEL (n)

```
for  $j = 0, \dots, k-1$  do
  for each  $s \in S$  do
     $\pi_j(s) = \operatorname{argmax}_{a \in A} \mathcal{Q}_j(s, a)$ 
    for each  $a \in A$ 
       $\hat{T} \mathcal{Q}(s, a) = r(s, a) + \gamma (\hat{P}^{\pi_j} \mathcal{Q}_j)(s, a)$ 
    end
     $\mathcal{Q}_{j+1}(s, a) = \hat{T} \mathcal{Q}_j(s, a)$ 
  end
end
end
Return  $\mathcal{Q}_k$ 
```

Model based policy iteration:

Input: reward, γ , π_0 , n , k .

$\hat{P} = \text{ESTIMATE_MODEL}(n)$

$$Q_0 = (I - \gamma \hat{P}^{\pi_0})^{-1} r;$$

for $j = 0, 1, \dots, k-1$

for each $s \in S$ do

$$\pi_j(s) = \underset{a \in A}{\operatorname{argmax}} Q_j(s, a)$$

end

$$\hat{Q}^{\pi_j} = (I - \gamma \hat{P}^{\pi_j})^{-1} r$$

$$Q_{j+1} = \hat{Q}^{\pi_j}$$

end

return Q_k ;

denote $S \times A$ by Z ;

MODELESTIMATE (a):

$\forall (s, z) \in S \times Z$ set $m(s, z) = 0$

for each $z \in Z$

for $i = 1, \dots, n$ do

 $s \sim P(\cdot | z)$

$m(s, z) = m(s, z) + 1$

end.

$\forall s \in S, \hat{P}(s | z) = \frac{m(s, z)}{n}$

end

return \hat{P}

ASSUMPTIONS: $Z = S \times A$ is finite; Assume

$$r_z(s, a) \in [0, 1];$$

THEOREM: \exists constants c, c_0, d & d_0 s.t.
 $\forall \epsilon \in (0, 1), \forall \delta \in (0, 1)$, a total sampling budget

$$T = \left\lceil \frac{c |S| |A|}{(1-\gamma)^3 \cdot \epsilon^2} \log \left(\frac{c_0 |S| |A|}{\delta} \right) \right\rceil$$

Suffices for $\|Q^* - Q_k\|_{\infty} \leq \epsilon$ with probability $1 - \delta$, after $k = \left\lceil \frac{d \log \left(\frac{d_0}{\epsilon(1-\gamma)} \right)}{\log \left(\frac{1}{\gamma} \right)} \right\rceil$ iterations of QVI or PI

The above theorem holds for finding a near optimal policy

The bounds are almost tight;

