

• Analysis of the Q_{max} algorithm:

Note that we only really need a π which tells us what to do for H steps - so it need not be a policy at all.

(ie) a stationary policy.

However what we show is that the optimal π_k^* obtained from M_k is good! So in fact we get a policy, a stationary policy.

- This π_k^* is what we use till k is updated.

- We could even work with an optimal H step policy in M_k . That too would work.

- But we work with the optimal policy in M_k .
- In fact with the "empirical" \hat{M}_k , and the optimal policy there since we are only estimating M_k .

LAST TIME:

Let M be an MDP, K -known states. M_k , the induced MDP. For any stationary policy π and state $s \in S$,

$$V_{M_k}^{\pi}(s) \geq V_M^{\pi}(s) \quad \&$$

$$V_M^{\pi}(s) \geq V_{M_k}^{\pi}(s) - \Pr_K[\text{escape from } s \text{ to } S]$$

CDR:

$$V_M^{\pi^*(M_K)}(s) \geq V_M^*(s) - P_{+1}^{\pi^*(M_K)} \left[\text{escape from } K \mid s_0 = s \right]$$

PF:

$$V_M^{\pi^*(M_K)}(s) \geq V_{M_K}^{\pi^*(M_K)}(s) - P_{\downarrow M}^{\pi^*(M_K)} \left[\text{escape from } \downarrow \mid s_0 = s \right]$$

$$\geq V_{M_K}^{\pi^*(M)}(s) - P_{\downarrow M}^{\pi^*(M_K)} \left[\text{escape from } K \mid s_0 = s \right]$$

$$\geq V_M^{\pi^*(M)}(s) - P_{\downarrow M}^{\pi^*(M_K)} \left[\text{escape from } K \mid s_0 = s \right]$$

- we will prove the main theorem.

Main theorem:

Let s_t be the state visited at round t , & let $m = O\left(\frac{H^2}{\epsilon^2} \log\left(\frac{S^2 A}{\delta}\right)\right)$. For any $\epsilon > 0$, $\delta < 1$, w.p

$1 - \delta$, $V_M^{\pi_t}(s_t) \geq V_M^*(s_t) - \epsilon$ for all but

$O\left(\frac{H^3 S^2 A}{\epsilon^3} \log\left(\frac{S^2 A}{\delta}\right)\right)$ rounds in the MDP.

Proof:

We showed that the value function

will be ϵ away from the optimal for at most

$$O\left(\frac{m H |S| |A|}{\epsilon} \log\left(\frac{|S| |A|}{\delta}\right)\right)$$

rounds.

- We do not have the Markov Chain.
- We only sample from it!

So we cannot really compute M_k , only an approximation \hat{M}_k .

And we will simulate that!

We showed if:

$$\sum_{s' \in S} |P_{\mu}(s'|s, a) - P_{\mu'}(s'|s, a)| \leq \epsilon, \quad \forall s, a \text{ \& } s' \in S$$

$$|q_{\mu}(s, a) - q_{\mu'}(s, a)| \leq \epsilon_1 \quad \forall s, a \text{ then}$$

for every policy π (stationary)

$$\forall s, \quad |V_{\mu}^{\pi}(s) - V_{\mu'}^{\pi}(s)| \leq \frac{\gamma}{1-\gamma} \epsilon_1 + \epsilon_2$$

• We'll use this, assuming that $\underline{\varepsilon}_1 = \varepsilon/2$.

Then: $\forall t$,

$$\left| V_{M_k}^{\pi_t}(s) - V_{M_k}^{\pi_t^A}(s) \right| \leq \varepsilon/2$$

$\therefore \forall s$,

$$\underbrace{V_{M_k}^{\pi_t}(s)} \geq V_{M_k}^{\pi_t^A}(s) - \frac{\varepsilon}{2} \geq \underbrace{V_{M_k}^{\pi^*}(s) - \frac{\varepsilon}{2}} \geq V_{M_k}^{\pi^*}(s) - \varepsilon$$

π_t^A is optimal w.r.t M_k^A

• So, $\forall s$,

$$V_M^{\pi_t}(s) \geq V_{M_k}^{\pi_t^A}(s) - \mathbb{P}_{\mathcal{G}_M}^{\pi_t^A}[\text{escape from } k \mid s_0 = s]$$

$$\geq V_{M_k}^{\pi^*}(s) - \mathbb{P}_{\mathcal{G}_M}^{\pi_t^A}[\text{escape from } k \mid s_0 = s] - \varepsilon$$

$$\geq V_M^{\pi^*}(s) - \varepsilon - \mathbb{P}_{\mathcal{G}_M}^{\pi_t^A}[\text{escape from } k \mid s_0 = s]$$

1. If $P_{\sigma, \mu}^{\text{TE}}[\text{escape from } K \mid s_0 = s] \leq \epsilon$, we are within 2ϵ of the optimal value.

O.W. we have a good chance of escaping!
and then we know for all best

$O\left(\frac{m K |S(A)|}{\epsilon} \log \frac{|S(A)|}{\delta}\right)$ rounds we

are within ϵ ,

\wedge we are within 2ϵ \forall but rounds.

Now to find m , so M_K & \hat{M}_K are close!

Lemma:

Assume m samples are obtained from a distribution p , whose support is of size N .
 \hat{p} , empirical distribution.

If $m = O\left(\frac{N^2}{\epsilon^2} \log\left(\frac{N}{\delta}\right)\right)$, with prob $1 - \delta$,

$$\sum_i |\hat{p}(i) - p(i)| \leq \epsilon.$$

* Can improve this. Will be in H.W.

We select m samples; And for each i we calculate,

$$\frac{\# \text{ i's seen}}{m};$$

If we ensure $\Pr\left[|\hat{p}(i) - p(i)| \leq \frac{\epsilon}{\sqrt{2}}\right] \geq 1 - \delta$

then with prob $(1 - \delta)$, $\forall i, \Pr\left[|\hat{p}(i) - p(i)| \leq \frac{\epsilon}{\sqrt{2}}\right] \geq 1 - \delta$

and so with prob $1-\delta$,

$$\sum_i \left| \hat{p}(i) - p(i) \right| \leq \frac{\epsilon}{N}, N \leq \epsilon.$$

For one i , use Hoeffding;

$$\Pr \left[\left| \sum_k X_k^i - \sum_k \mathbb{E}(X_k^i) \right| \geq t \right] \leq \exp \left(\frac{-2t^2}{\sum_{k=1}^m (b_k^i - a_k^i)^2} \right)$$

Here $b_k^i - a_k^i = 1$.

So: $\Pr \left[\left| \frac{\sum X_k^i}{m} - p_i \right| \geq \frac{t}{m} \right] \leq \exp \left(\frac{-2t^2}{m} \right)$

$$\frac{t}{m} = \frac{\epsilon}{N} \quad \therefore t = \frac{\epsilon m}{N};$$

$$\leq \exp \left(- \frac{2 \cdot 4 \cdot \epsilon^2 m^2}{N^2 \cdot m} \right)$$

$$\leq \frac{1}{\exp \left(\frac{8\epsilon^2 m}{N^2} \right)} \leq \frac{\delta}{2} \quad \text{if}$$

$$m \geq \frac{N^2}{8\epsilon^2} \log \left(\frac{2}{\delta} \right)$$

Want:

$$\sum_{s'} \left| P_{M_{12}}(s'|s, a) - P_{\hat{M}_{12}}(s'|s, a) \right| \leq \epsilon,$$

then: $\|V_{M_{12}}^{\pi_H}(s) - V_{\hat{M}_{12}}^{\pi_H}(s)\| \leq \underbrace{\frac{\gamma}{1-\gamma}}_{\leq \epsilon/2} \epsilon,$

Want:

What about ϵ_1 ?

We want for all consecutive transitions that the two Markov chains agree to ϵ_1 of their transition probabilities.

But error accumulates to $\epsilon_1 H$ in H rounds.

So we can only guarantee $\frac{\gamma}{1-\gamma} \epsilon_1 H$!

Cloeness!

\therefore set $\frac{\gamma}{1-\gamma} \epsilon_1 H = \frac{\epsilon}{2}$

$$\therefore \boxed{\epsilon_1 = \frac{1-\gamma}{\gamma} \frac{\epsilon}{2H}}$$

δ

Want this confidence for all (s, a)

\therefore We can afford an $\frac{\delta}{|S||A|}$ error for one (s, a) , and this is what we plug into

$$m \geq \frac{N^2}{8 \epsilon_0^2} \log \left(\frac{N}{\delta_0} \right)$$

Finally:

$$\epsilon_0 = \frac{1-\gamma}{\delta} \frac{\epsilon}{2H} ; \delta_0 = \frac{\delta}{|S||A|}$$

$$\therefore m \geq \frac{N^2 \cdot 4H^2}{8 \epsilon^2} \left(\frac{\delta^2}{(1-\gamma)^2} \right) \log \left[\frac{N |S||A|}{\delta} \right]$$

Here $N = |S||A|$.

$$\therefore \left(\frac{S^2 A^2 H^2}{\epsilon^2} \log \left[\frac{|S||A|^2}{\delta} \right] \right)$$

depends to:

$$O\left(\frac{|S(A)|^L}{\epsilon^L} \log \frac{|S(A)|}{\delta}\right).$$

- Next chapter is policy gradient.
- We look at more Monte Carlo methods.
See Barto & Sutton.

9 ESTIMATING STATE-VALUE FUNCTION for a policy π (* π is fixed)

Start from $s \rightarrow$

Generate a history $(s_0=s, A_0, R_0, s_1, a_1, r_1, \dots)$

Compute $\sum_{k=0}^{\infty} \gamma^k R_k$

We are constructing an "unbiased" estimator
of $v^\pi(s)$

Alg:

Input π :

$v \leftarrow 0$; $R(s) \leftarrow \text{null } \forall s$;

for $i = 1, \dots, k$

Generate an episode $s_{0,i}, a_{0,i}, r_{0,i}, s_{1,i}, a_{1,i}, r_{1,i}, \dots$

for each state s appearing in episode,
 $t \leftarrow$ time of first occurrence of s .

$$G_{i,s} \leftarrow \sum \gamma^k R_{t+k}$$

Append G_i to $R(s)$

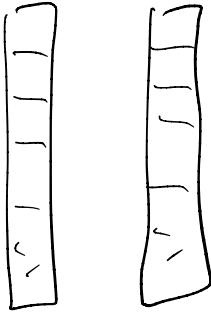
end

$$v(s) = \frac{1}{k} \sum G_i$$

Q: Can we compute this using last occurrence of s ?

• $\mathbb{E}(G_i | s) = v^\pi(s);$

• i-i-d samples $\therefore \mathbb{E}\left(\frac{\sum G_i}{n}\right) = v^\pi(s);$



finite states \therefore finite vectors;

Works if every state is visited infinitely often

A variant:

Every visit Monte Carlo:

- for each state s appearing in the episode and each time t it occurs,

- find the ^{discounted} reward from that time
- append all such rewards to the list (s) ;
- Average of the rewards is $\hat{V}_\pi(s)$;

Gradient based Monte Carlo:

- Start with $v \in \mathbb{R}^{|S|}$; At time t

$$v(s_t) \leftarrow v(s_t) + \alpha [G_t - V(s_t)]$$
- Mean square error $(v) = \frac{1}{2} \mathbb{E}_S [(v^\pi(s) - v(s))^2]$

Gradient descent:

$$\begin{aligned}
 v &\leftarrow v - \alpha \frac{\partial \text{MSE}(v)}{\partial v} \\
 &= v - \alpha \mathbb{E} \left((v^\pi(s) - v(s)) (-1) \frac{\partial v(s)}{\partial v} \right)
 \end{aligned}$$

$$= v + \alpha \mathbb{E} \left[\left(v^\pi(s) - v(s) \right) \frac{\partial v(s)}{\partial v} \right]$$

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix} \leftarrow s.$$

$$V_{\text{new}}(s_t) = V_{\text{old}}(s_t) + \alpha \left(G_t - V_{\text{old}}(s_t) \right)$$

$$\mathbb{E}(G_t | s_t) = v^\pi(s_t);$$