

---

---

---

---

---



Recall  $k$ -max.  $\gamma$  algorithm.

•  $\hat{M}_k$  - induced MDP:

Def: Let  $M = (S, A, P, r, \gamma)$ ;  $K \subseteq S$  - known.

$$\forall s \in K \quad P_{M_k}(s' | s, a) = P_M(s' | s, a) \quad \& \quad r_{M_k}(s, a) = r_M(s, a)$$

$\forall s \notin K$ :

$$P_{M_k}(s' | s, a) = \mathbb{1}(s' = s) \quad \& \quad r_{M_k}(z | s, a) = \mathbb{1}(z = 1)$$

↑  
Stay back in the same  
state

↑  
Reward = 1.

$$\text{Set } H = \frac{\log\left(\frac{2}{\epsilon(1-\gamma)}\right)}{1-\gamma}$$

↑  
WAS THE # iterations for  $Q$ -value-iteration  
to converge with  $\sqrt{T^k}$  with  $\epsilon$  ball of  $V^*$  (in  $\ell_\infty$  norm)

Thm: Let  $s_t$  be the state visited by  $\epsilon$ -max- $\delta$  algorithm in round  $t$  and set  $m = \left( \frac{SA^2}{\epsilon^2} \log\left(\frac{SA}{\delta}\right) \right)$

for  $\epsilon > 0, \delta < 1$ ,

$$V_M^{\pi^*}(\epsilon) \geq V_M^*(s_t) - \epsilon \quad \forall \text{ but } O\left(\frac{SA^3}{\epsilon^3} \log\left(\frac{SA}{\delta}\right)\right)$$

rounds in the MDP.

← We're only talking of  $V(s_t)$  in round  $t$ .

- We may want to prove that it finds a near-opt policy, (ie) a policy whose expected reward is within  $\epsilon$  of the optimal, when taking expectations over the start state as well.

↑  
WE DON'T PROVIDE SUCH A GUARANTEE!

- OPTIMISM in the face of uncertainty.

\*: In the proof below  $r$  also refers to expected reward!

Lemma:

Let  $M, M'$  be MDPs on same action & state spaces. If

$$\sum_{s' \in S} |P_M(s'|s, a) - P_{M'}(s'|s, a)| \leq \epsilon_1, \quad \forall s, a.$$

&

$$|r_M(s, a) - r_{M'}(s, a)| \leq \epsilon_2 \quad \forall s, a$$

then

$$\forall \pi, \text{ stationary policies, } \|V_M^\pi - V_{M'}^\pi\|_\infty \leq \frac{\gamma \epsilon_1}{1 - \gamma} + \epsilon_2$$

Proof:

$$\begin{aligned} |V_M^\pi(s) - V_{M'}^\pi(s)| &\leq \\ & (1 - \gamma)\epsilon_2 + \gamma \left| \sum_{s' \in S} P_M(s'|s, \pi(s)) V_M^\pi(s') - P_{M'}(s'|s, \pi(s)) V_{M'}^\pi(s') \right| \\ & \downarrow \\ & \leq (1 - \gamma)\epsilon_2 + \gamma \left( \sum_{s' \in S} P_M(s'|s, \pi(s)) [V_M^\pi(s') - V_{M'}^\pi(s')] \right. \\ & \quad \left. + \sum_{s' \in S} [P_M(s'|s, \pi(s)) - P_{M'}(s'|s, \pi(s))] V_{M'}^\pi(s') \right) \end{aligned}$$

$$\leq (1-\gamma) \epsilon_2 + \gamma \|V_M^\pi(s) - V_{M^1}^\pi(s)\|_\infty + \gamma \epsilon_1$$

$$\therefore \|V_M^\pi(s) - V_{M^1}^\pi(s)\|_\infty \leq \epsilon_2 + \frac{\gamma}{1-\gamma} \epsilon_1$$

==>

NOTATION:

$$\overline{P_M^\pi}[\text{escape from } K | s_0 = s] = \mathbb{1}(s \notin K) + \sum_{t=1}^{\infty} \gamma^t P_M^\pi(s_t \notin K, s_0, \dots, s_{t-1} \in K)$$

discounted

The above term is the discounted probability of reaching an unknown state when executing  $\pi$  starting from state  $s$ .

### Lemma 3.4:

Let  $M$  be an MDP;  $K$ -known state &  $M_K$  as before.  $\forall \pi$ , stationary &  $s \in S$ ,

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

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

- $M_K$  is an optimistic version of  $M$  since it gives greater value to all states  $\forall$  stationary policies.
- But **value of optimistic policy not too high!** The difference is high only if there is a large probability of escaping to an unknown state!

•  $V_{M_K}^\pi(s) \geq V_M^\pi(s)$  is clear

• If  $s \notin K$ , get a max reward of 1  
 $(1-\gamma)(1 + \gamma + \gamma^2 + \dots) = 1$

• If  $s \in K$  - immediate reward is exactly that in  $M$   
 - subsequent steps - remain in  $K$   
 and get same reward  
 - leaving  $K$  &  
 get maximum reward!

•  $\left| V_M^\pi(s) - V_{M_K}^\pi(s) \right| \leq$

$$\mathbb{1}(s \notin K) + \mathbb{1}(s \in K) \gamma \left[ \sum_{s' \in S} P_M(s' | s, \pi(s)) V_M^\pi(s') - P_{M_K}(s' | s, \pi(s)) V_{M_K}^\pi(s') \right]$$

Use  $\because s \in K, P_M(s' | s, \pi(s)) = P_{M_K}(s' | s, \pi(s))$

$$\leq \mathbb{1}(s \notin K) + \mathbb{1}(s \in K) \gamma \left[ \sum_{s'} P_M(s' | s, \pi(s)) \left( \frac{V_M^\pi(s') - V_{M_K}^\pi(s')}{V_M^\pi(s')} \right) \right]$$

$$\leq \mathbb{1}(s \notin K) + \mathbb{1}(s \in K) \sigma P_M(s' \notin K | s, \pi(s))$$

$$+ \mathbb{1}(s \in K) \sigma \left| \sum_{s' \in K} P_M(s' | s, \pi(s)) \left( \underbrace{V_M^\pi(s') - V_{M_K}^\pi(s')} \right) \right|$$

expand this again

- s' \in K, get

$$\sigma \sum_{s'' \in S} \mathbb{1}(s' \in K) \left[ P_M(s'' | s', \pi(s')) V_M^\pi(s'') - P_{M_K}(s'' | s', \pi(s')) V_{M_K}^\pi(s'') \right]$$

which:

if  $s'' \notin K$  gives  $\sigma P_M(s'' \notin K | s', \pi(s'))$

$\hookrightarrow s \in K$

if  $s'' \in K$  gives

$$\sigma \sum_{s''' \in S} P_M(s''' | s'', \pi(s'')) \left[ V_M^\pi(s''') - V_{M_K}^\pi(s''') \right]$$

Coeff of  $\sigma^2$ ,

$$\sigma^2 \mathbb{1}(s \in K) \mathbb{1}(s' \in K) \left[ \text{Prob of reading } s'' \notin K \text{ in 2 steps} \right]$$



• (1c) Prob of escaping from  $K$  in 2 steps  $\approx \delta^2$ .

• As we expand this:

$$\leq \mathbb{1}(s \notin K) + \mathbb{1}(s \in K) \cdot \delta P_{\pi} (s' \notin K \mid s, \pi(s))$$

$$+ \mathbb{1}(s \in K) \mathbb{1}(s' \in K) \delta^2 P_{\pi} (s'' \notin K \mid s, s' \in K)$$

$$+ \mathbb{1}(s \in K) \mathbb{1}(s' \in K) \mathbb{1}(s'' \in K) \delta^3 P_{\pi} (s''' \notin K \mid s, s', s'' \in K)$$



exactly  $P_{\pi}^{\Pi} [\text{escape from } K \mid s_0 = s]$

\* Jump

- Lemma

- With sufficiently many visits to states with a large escape probability, all states become known. ) whp
- With high probability the # rounds following a visit where the policy's value function is significantly suboptimal is at most  $H$ .

COR:

Implicit Explore - Exploit:

$$V_M(s) \geq V_M^*(s) - P_M \left[ \text{escape for } K \mid s_0 = s \right]$$

From lemma:

$$\overline{V_M^{\pi^*(M_k)}} \geq V_{M_k}^*(s) - P_M \left[ \text{escape for } K \mid s_0 = s \right]$$

Why  $\rightarrow$   $\geq V_{M_k}^{\pi^*(M)}(s) - P_M \left[ \text{escape for } K \mid s_0 = s \right]$

Why  $\rightarrow$   $\geq V_M^{\pi^*(M)}(s) - P_M \left[ \text{escape for } K \mid s_0 = s \right]$

Policy computed in each episode is near optimal

## Lemma:

with probability at least  $1 - \delta$ , the # rounds  $t$  with  $V_M^{\pi_t}(s) \leq V_M^*(s) - \epsilon$  is at most  $O\left(\frac{mHSA}{\epsilon} \log\left(\frac{1}{\delta}\right)\right)$

## Proof:

2 steps.

1) If in round  $t$ , the prob. of escape from  $K$  starting from  $s_t$  is large, then we have a large escape probability in the next  $H$  steps.

2) In fact, with high probability, the algorithm encounters a unknown state in the next  $H$  steps.

i)  $s = s_t, \pi_t$ ; <sup>Suppose</sup>  $P_M^{\pi_t}[\text{escape from } K | s = s] \geq \epsilon$

Define:  $p_H = \mathbb{1}(s \notin K) + \sum_{t=1}^H P_M^{\pi_t}(s_t \notin K | s_0, \dots, s_{t-1} \in K)$

Prob of escaping from  $K$  in  $t$  steps

- Let  $s_t$  = state encountered at round  $t$   
&  $\pi_t$  be the policy.

Suppose  $\mathbb{P}_{\pi}^{\pi} [\text{escape from } K \mid s_0] \geq \epsilon$ .

- Set  $p_H = \mathbb{1}(s \notin K) + \sum_{t=1}^H \mathbb{P}_{\pi}^{\pi} (s_t \notin K \mid s_0, \dots, s_{t-1})$

Probability of escaping from  $K$  in  $H$  steps!

No discounted probability!

$$\epsilon \leq \mathbb{1}(s \notin K) + \sum_{t=1}^{\infty} \gamma^t \mathbb{P}_{\pi}^{\pi} (s_t \notin K \mid s_0, \dots, s_{t-1})$$

$$\leq p_H + \sum_{t=H+1}^{\infty} \gamma^t \leq p_H + \frac{\gamma^{H+1}}{1-\gamma}$$

for  $H$  as chosen  $\frac{\gamma^{H+1}}{1-\gamma} \leq \frac{\epsilon}{2}$

for this  $H$   $p_H \geq \frac{\epsilon}{2}$

Bound

# actions before we have enough visits  
to unknown states.

=

Define  $t_1, t_2, \dots$  rounds s.t

$$\& \quad |t_{i+1} - t_i| \geq H.$$

If  $\pi_i$  - policy used in time  $i$ ,  $K_i$  - <sup>known</sup> states,

then

$$\mathbb{P}^{\pi_i} [\text{escape from } K_i \mid s_0 = s_{t_i}] \geq \epsilon.$$

=

Set:

$$X_i = \mathbb{1} \left( \exists s \text{ in } \{s_{t_i}, s_{t_i+1}, \dots, s_{t_i+H}\} : s \notin K_i \right)$$

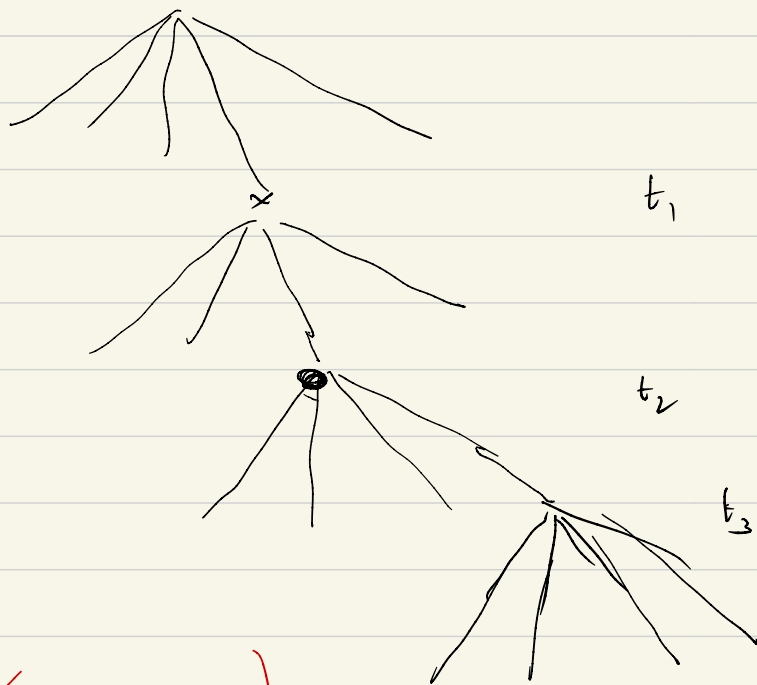
=

By def:  $\mathbb{E} (X_i \mid s_{t_i}) \geq \epsilon.$

Let  $\mathcal{F}_i$  - values of all random variables prior to time  $t_i$ , including time  $t_i$

Clearly  $\mathcal{F}_0 \subseteq \mathcal{F}_1 \subseteq \dots$  a  $\sigma$ -field;

▣  $X_i$  is measurable w.r.t  $\mathcal{F}_i$



What  $E(X_i | \mathcal{F}_{i-1})$

At time  $t_i$  — use policy  $\pi_i$ .

And  $x_i = 1$  iff we escape from  $K_i$  in  $H$  steps; conditioned on  $\mathcal{F}_{i-1}$ .

$$- \mathbb{E} [x_i | \mathcal{F}_{i-1}] \geq \frac{\epsilon}{2};$$

Now

$$\begin{aligned} & \mathbb{E} \left[ (x_i - \mathbb{E}[x_i | \mathcal{F}_{i-1}])^2 \mid \mathcal{F}_{i-1} \right] \\ &= \mathbb{E} \left[ x_i^2 - 2x_i \mathbb{E}[x_i | \mathcal{F}_{i-1}] + \mathbb{E}[x_i | \mathcal{F}_{i-1}]^2 \mid \mathcal{F}_{i-1} \right] \\ &= \mathbb{E} [x_i^2 | \mathcal{F}_{i-1}] - \mathbb{E} [x_i | \mathcal{F}_{i-1}]^2 \\ &\leq \mathbb{E} [x_i^2 | \mathcal{F}_{i-1}] = \mathbb{E} [x_i | \mathcal{F}_{i-1}] = x_i = \begin{cases} 0 \\ 1 \end{cases} \end{aligned}$$

## Freedman's inequality:

Let  $X_1, X_2, \dots, X_T$  be a sequence of real valued random variables adapted to the filtration  $\mathcal{F}_t$ .  
 $\therefore X_i$  is measurable w.r.t  $\mathcal{F}_i$  & further assume that  $\mathbb{E}[X_i | \mathcal{F}_{i-1}] < \infty$ ;

Define  $S = \sum_{t=1}^T X_t$ ,  $V = \sum_{t=1}^T \mathbb{E}(X_t^2 | \mathcal{F}_{t-1})$  and let

$X_t \leq R$  almost surely  $\forall t$ ;

$\forall \delta \in (0, 1)$  &  $\lambda \in [0, 1/2]$  with  $\lambda \geq \delta$  at least  $1-\delta$ ,

$$S \leq (e-2)\lambda V + \frac{\ln(1/\delta)}{\lambda}$$

Choosing:  $\lambda = \min\left(\frac{1}{2}, \sqrt{\frac{\ln(1/\delta)}{V}}\right)$  we get

$$S \leq 2\sqrt{V \ln(1/\delta)} + R \ln(1/\delta).$$



• Set  $\gamma_i = \mathbb{E}[x_i | \mathcal{F}_{i-1}] - x_i$

Applying Friedman's inequality:

$$\sum_{i=1}^n \gamma_i$$

$$= \sum_{i=1}^n \mathbb{E}[x_i | \mathcal{F}_{i-1}] - x_i$$

$$\leq 2 \sqrt{\ln\left(\frac{1}{\delta}\right) \sum_{i=1}^n \mathbb{E}\left[\mathbb{E}\left[(x_i | \mathcal{F}_{i-1}) - x_i\right]^2 \mid \mathcal{F}_{i-1}\right]} + \ln\left(\frac{1}{\delta}\right)$$

$$\leq 2 \sqrt{\ln\left(\frac{1}{\delta}\right) \sum_{i=1}^n \mathbb{E}[x_i^2 | \mathcal{F}_{i-1}]} + \ln\left(\frac{1}{\delta}\right)$$

$$\leq \frac{1}{2} \sum \mathbb{E}[x_i^2 | \mathcal{F}_{i-1}] + 3 \ln\left(\frac{1}{\delta}\right)$$

$$\circ \circ \quad \sum x_i \geq \frac{1}{2} \sum_{i=1}^n \mathbb{E}[x_i^2 | \mathcal{F}_{i-1}] - 3 \ln\left(\frac{1}{\delta}\right)$$

$$\geq \frac{n\epsilon}{2} \cdot \frac{1}{2} - 3 \ln\left(\frac{1}{\delta}\right)$$

Want:  $\sum_{i=1}^n x_i \geq mSA$

$$\therefore \frac{n\epsilon}{4} - 3 \ln\left(\frac{1}{\delta}\right) \geq mSA$$

$$n \geq \frac{4}{\epsilon} \left[ mSA + 3 \ln\left(\frac{1}{\delta}\right) \right]$$


---



---

• for rounds  $t \in [t_{left}, t_{right} - 1]$  - prob of escape  $\leq \epsilon$

∴ value function  $\bar{c}_n^{\text{max}}$  optimal on these rounds.