



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:

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

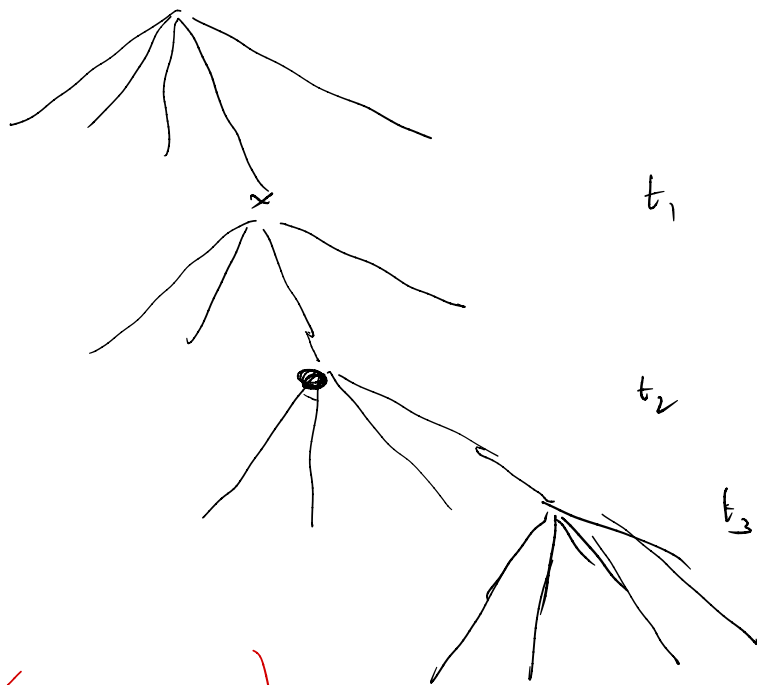
=

By def:  $\mathbb{E} \left( X_i \mid s_{t_i} \right) \geq \frac{\epsilon}{2}$

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

Clearly  $\mathcal{F}_\emptyset \subseteq \mathcal{F}_1 \subseteq \dots$  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  $\mathcal{F}_{i-1}$ .

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

Now

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

## 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  $\chi_i = \mathbb{E}[x_i | \mathcal{F}_{i-1}] - x_i$

Applying Freedman's inequality:

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

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

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

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

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

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

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

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

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

$$n \geq \frac{4}{\varepsilon} \left[ mSA + 3 \ln\left(\frac{1}{\delta}\right) \right] = \text{Set } n = \frac{4(mSA + 3 \ln(1/\delta))}{\varepsilon}$$


---

- for rounds  $t \in [t_{\text{left}}, t_{\text{right}} - 1]$  - prob of escape  $\leq \varepsilon$   
 $\therefore$  value function  $\hat{c}_n^{\text{max}}$  optimal on those rounds.

So with prob  $(1-\delta)$ ,  $\sum x_i \geq mSA$  which means all states will be known.

• The # states which are not  $\epsilon$  optimal?

Each of the rounds  $t_i$ ,  $i=1, \dots, n$

Gives us  $O\left(\frac{m|S||A|}{\epsilon} \ln(|S|)\right)$  states.

And for each of them  $H$  more rounds.

$\therefore$  At most  $O\left(\frac{m|S||A|}{\epsilon} H \ln(|S|)\right)$  states.

-Note: This is independent of the # episodes.

Where an episode is that period where  $K$  remain the same.



## Weaker bound

EXPLAINING ALL THIS WITHOUT FRIEDMAN.

^  
• What we have shown is that  
the # rounds where  $V_M^{\pi_t}(s_t) \leq V_M^{\pi^*}(s_t) - \epsilon$

is at most  $O\left(\frac{mHSA}{\epsilon} \ln(1/\delta)\right)$ .

Note: Think of an episode as

the rounds where  $K$  remains the same.

(1c) policy does not change.

For the value  $H$  chosen, over  $t$ 's are such  
that  $\text{scope}(s_t) = 1$  is at least  
 $\epsilon/2$  [  $x_i = 1$  ]

- We want to estimate # episodes where  
we may fail;

# state action pairs, =  $|S| |A|$ .  
and a state is known if we see each

each state action pair  $m$  times.

$\therefore$  An upper bound on # state action pairs is  $m|S||A|$ .

If we have  $\frac{m|S||A|}{\epsilon}$  episodes, we expect more than  $\frac{m|S||A|}{\epsilon} \times \epsilon = m|S||A|$  successes. In which

case we are done

$$\text{Set } n = \frac{m|S||A|}{\epsilon} \log\left(\frac{|S||A|}{\delta}\right)$$

Then <sup>the</sup> expected # successes  $\geq \frac{m|S||A|}{\epsilon} \log\left(\frac{|S||A|}{\delta}\right)$

- All  $m|S||A|$  explorations succeed with probability  $\geq 1 - \delta$ ;

Main theorem:

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

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

$O\left(\frac{1}{\epsilon^3} \frac{S^2 A^2}{\delta} \log\left(\frac{S^2 A}{\delta}\right)\right)$  rounds in the MDP.

Proof:

Assume  $m$  is large.

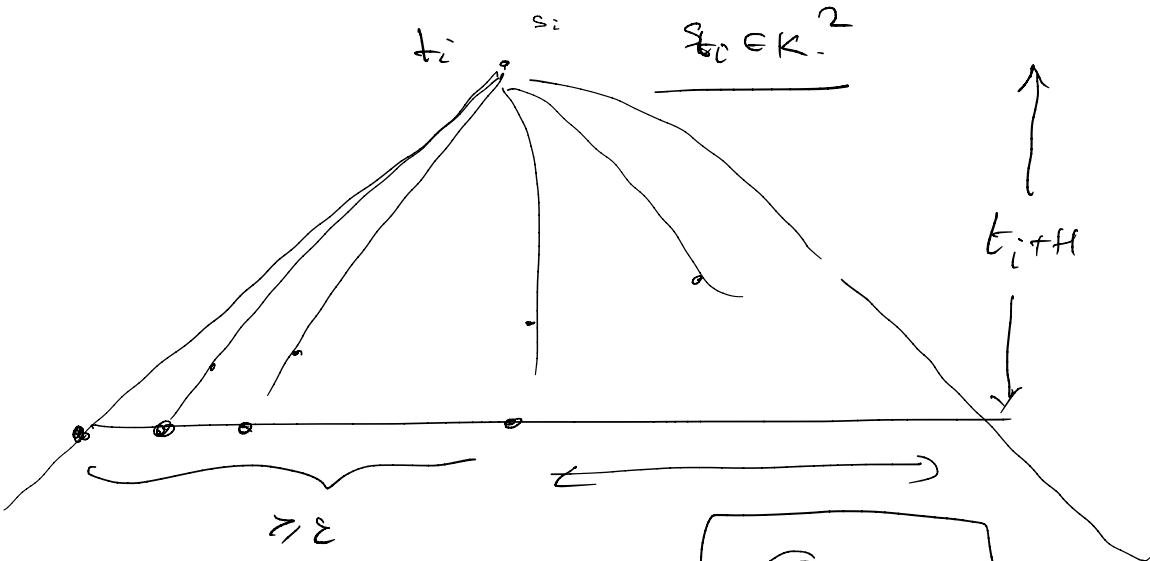
And we have an approximation  $\hat{M}_K$  of the  
includ draw.

Assume that  $\sum |P_{M_K}(s' | s, a) - P_{\hat{M}_K}(s' | s, a)| \leq \epsilon$ ,

so that  $\|V_{M_K}^{\pi} - V_{\hat{M}_K}^{\pi}\| \leq \frac{\delta}{1-\delta} \epsilon = \frac{\epsilon}{2}$ .

$$(x)_{t \leq t_i - 1}$$

value of all random variable.



$$V_{M,K}^{\pi_t}(s) \geq V_{\hat{M}_K}^{\pi_t}(s) - \frac{\epsilon}{2}$$

$\epsilon/2$  bound

$\frac{1}{2}$   
# number of steps

Visits unknown state with prob  $\geq \frac{\epsilon}{2}$

$$\epsilon \leq P[\sigma \in R] + \gamma P[\sigma \in R, S \in K] + \gamma^2 P - -$$