

March 18) Batch RL II.

① Better error control for model misspec.

$$M = (P, r) \quad \hat{M} = (\hat{P}, \hat{r})$$

$\downarrow \hat{\pi}$

$$H = \frac{1}{1-\delta}$$

$$\sqrt{\hat{\pi}} \geq \sqrt{\pi^*} - \delta \sqrt{1}$$

$$\delta \leq \frac{1}{(1-\delta)^3} \left(\|P - \hat{P}\| + \|r - \hat{r}\| \right)$$

$$\Rightarrow n \approx \frac{H^6 SA}{\varepsilon^2} \log\left(\frac{1}{\delta}\right) \quad \frac{\sqrt{1/n}}{HSA} \leq \varepsilon$$

(s,a) - sampling

$$\frac{H^3 SA}{\varepsilon^2} \log\left(\frac{1}{\delta}\right)$$

2.

lower bound

3.

Data collection by following policies

Lower bound $S(A^H)$ $\xrightarrow{H^*}$ trajectories
 $\xrightarrow{A^{\pi_b}}$ to get $(\varepsilon, \frac{1}{4})$

$$S = H$$

(s,a)-sampling:

$$\frac{H^3 HA \log\left(\frac{1}{\delta}\right)}{\varepsilon^2}$$

days free labeled
of traj. adaptivity.

$$H = H_{\delta, \varepsilon}$$



$$H \cdot A^H$$

1/3rd of it.

behavior policy

(s,a) sample

planning

$$① M = (P, r)$$

$$\widehat{M} = (\widehat{P}, \widehat{r}) : \widehat{\pi}$$

$$\widehat{V}^{\widehat{\pi}} \approx \widehat{V}^*$$

$$\textcircled{A} V^* - V^{\widehat{\pi}} = V^* - \widehat{V}^* + \widehat{V}^* - \widehat{V}^{\widehat{\pi}} + \widehat{V}^{\widehat{\pi}} - V^{\widehat{\pi}}$$

$$\pi^* \text{ opt. in } M \rightarrow \leq \underbrace{V^* - \widehat{V}^*}_{\text{Opt. error}} + \underbrace{\widehat{V}^* - \widehat{V}^{\widehat{\pi}}}_{\text{Opt. error}} + \underbrace{\widehat{V}^{\widehat{\pi}} - V^{\widehat{\pi}}}_{\text{Opt. error}}$$

$$\pi \in \{\pi^*, \widehat{\pi}\}$$

$$\text{For simplicity: } r = \widehat{r}$$

$$\begin{aligned} \textcircled{A} \underbrace{V^{\pi} - \widehat{V}^{\pi}}_{\substack{\frac{1}{1-x} - \frac{1}{1-y} \\ = \frac{1-y-(1-x)}{(1-x)(1-y)} \\ = \frac{x-y}{(1-x)(1-y)}}} &= (I - \gamma P_{\pi})^{-1} r_{\pi} - (I - \gamma \widehat{P}_{\pi})^{-1} r_{\pi} \\ &= (I - \gamma P_{\pi})^{-1} \left[(I - \gamma \widehat{P}_{\pi}) - (I - \gamma P_{\pi}) \right] (I - \gamma \widehat{P}_{\pi})^{-1} r_{\pi} \\ &= \gamma (I - \gamma P_{\pi})^{-1} (P_{\pi} - \widehat{P}_{\pi}) (I - \gamma \widehat{P}_{\pi})^{-1} r_{\pi} \\ &= \gamma (I - \gamma P_{\pi})^{-1} M_{\pi} (P - \widehat{P}) \widehat{V}^{\pi} \end{aligned}$$

$$\|V^{\pi} - \widehat{V}^{\pi}\|_{\infty} \leq \frac{\gamma}{1-\gamma} \| (P - \widehat{P}) \widehat{V}^{\pi} \|_{\infty}$$

wp 1%:

$$\| (P - \widehat{P}) V^{\pi} \|_{\infty} \sim \sqrt{\frac{\log(SA/\zeta)}{n}} \frac{1}{1-\gamma} \frac{1}{1-\gamma}$$

Hoeffding

Loose!

Cheat!

$$V^* - V^{\hat{F}} \leq \frac{2\gamma}{(1-\gamma)^2} \sqrt{\frac{\log(\dots)}{n}} + \epsilon_{\text{opt}}$$

$\epsilon_{\text{opt}} \approx 0$

$$n \geq \frac{H^4 \log(\dots)}{\epsilon^2} \leq \epsilon$$

$H^G \rightarrow H^A$

$H^4 \rightarrow H^3 ?$

Hoeffding : worst-case tight!

Bernstein's \leq

$$X_1, \dots, X_n \in [0, b] , b > 0$$

i.i.d.

$$\bar{X}_n = \frac{1}{n} (X_1 + \dots + X_n) , \text{ w.p.t.s}$$

$$|\bar{X}_n - \mathbb{E}[X_1]| \leq \underbrace{5 \sqrt{\frac{2 \log(2/\delta)}{n}}}_{\text{red line}} + \underbrace{\frac{2b \log(2/\delta)}{3n}}_{\text{red line}}$$

$G^2 = \text{Var}(X_1)$

$$n \geq \frac{b^2 \log(\frac{1}{\delta})}{\epsilon^2}$$

$$b \sqrt{\frac{\log(2/\delta)}{2n}} \leq \epsilon$$

$$\text{Bemerkung: } n \geq \frac{\frac{6^2}{\varepsilon^2} \log(\dots)}{\varepsilon^2} \quad \leftarrow \frac{H^3}{\varepsilon^2}$$

$$n \geq \frac{b}{\varepsilon} \quad \frac{6^2}{\varepsilon^2} \geq \frac{b}{\varepsilon}$$

$$\delta^2 \leq b^2$$

$$1 = \frac{\sqrt{1-\gamma}}{\sqrt{1-\gamma}} =$$

$$\boxed{\frac{6^2}{b} \geq \varepsilon}$$

$$(P - \hat{P}) v^\pi$$

$$S'_1, \dots, S'_n \sim P_\alpha(s)$$

$$(P_\alpha(s) - \hat{P}_\alpha(s)) v^\pi = \underbrace{P_\alpha(s) v^\pi}_{\text{ }} - \frac{1}{n} \sum_{i=1}^n v^\pi(S'_i)$$

$$\hat{P}_\alpha(s, s') = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(S'_i = s')$$

$$\boxed{X_i = v^\pi(S'_i)}$$

$$\mathbb{E}[X_i] = P_\alpha(s)v^\pi$$

$$\text{Var}(X_i) = \sigma_\pi^2(s, a)$$

$$\sim \frac{1}{(1-\gamma)^2}$$

$$\boxed{\| (I - \gamma P_\pi)^{-1} M_\pi G_\pi \| \leq \sqrt{\frac{2}{(1-\gamma)^3}} = \frac{1}{1-\gamma} \frac{1}{\sqrt{1-\gamma}}}$$

discounted std. of value

$$\frac{1}{1-\gamma}$$

$$\sigma_\pi^2 \sim \frac{1}{1-\gamma}$$

2.

Can we improve upon

$$H^3(SA) \log\left(\frac{SA}{\epsilon}\right)$$

ϵ -subopt.

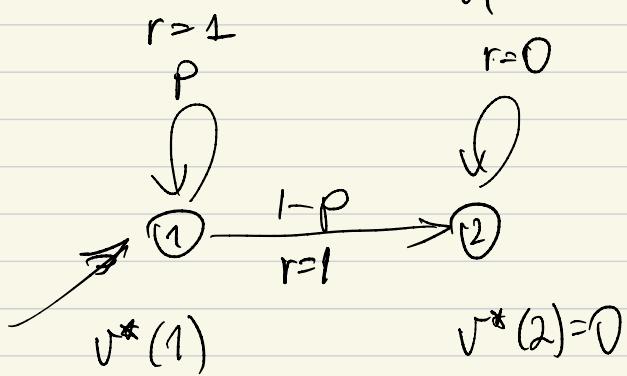
Sketch: $X_1, \dots, X_n \sim \text{Ber}(p)$ iid.

$$\rightarrow |\bar{X}_n - p| \approx \sqrt{\frac{p(1-p) \log(\frac{2}{\epsilon})}{2n}} + \dots$$

$\uparrow \quad \uparrow$

$$\text{Var}(X_i) = p(1-p)$$

$$\sqrt{\frac{p(1-p) \log(\frac{2}{\epsilon})}{2n}} \leq \epsilon \iff n \geq \frac{p(1-p) \log(\frac{2}{\epsilon})}{\epsilon^2}$$



$$V^*(1) = p(1 + \gamma V^*(1))$$

$$+ (1-p)(1 + \gamma \underbrace{V^*(2)}_{=0})$$

$$V_p^*(1) = 1 + p\gamma V^*(1)$$

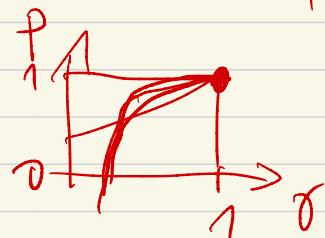
$$V_p^*(1) = \frac{1}{1-p\gamma}$$

$$\frac{d}{dp} V_p^*(1) = \frac{d}{dp} \frac{1}{1-p\gamma} = \frac{-(-\gamma)}{(1-\gamma p)^2} = \frac{\gamma}{(1-\gamma p)^2}$$

$$\epsilon \geq |V_P^*(1) - V_{\bar{X}_n}^*(1)| \approx \frac{d}{dp} V_P^*(1) \underbrace{|p - \bar{X}_n|}_{\substack{\text{=} \\ \text{---}}} = \frac{\tau}{(1-\tau p)^2} \sqrt{\frac{p(1-p)}{2n} \log(\dots)}$$

$$\tau \rightarrow 1$$

$$1 - \tau$$



$$1 - \tau p = 1 - \tau \left(1 + \frac{\tau-1}{2}\right) = 1 - \tau + \frac{\tau(\tau-1)}{2}$$

$$P = \frac{1+\tau}{2} = \frac{1}{2} + \frac{\tau}{2} = 1 + \frac{\tau-1}{2}$$

$$\underset{\sim}{\text{PDF}} \quad \frac{4\tau-1}{3\tau} = \frac{4}{3} - \frac{1}{3\tau}$$

$$= 1 + \frac{1}{3} \left(1 - \frac{1}{\tau}\right) = 1 + \frac{1}{3} \frac{1-\tau}{\tau}$$

$$1 - \tau p = 1 - \tau \left(1 + \frac{1}{3} \frac{1-\tau}{\tau}\right)$$

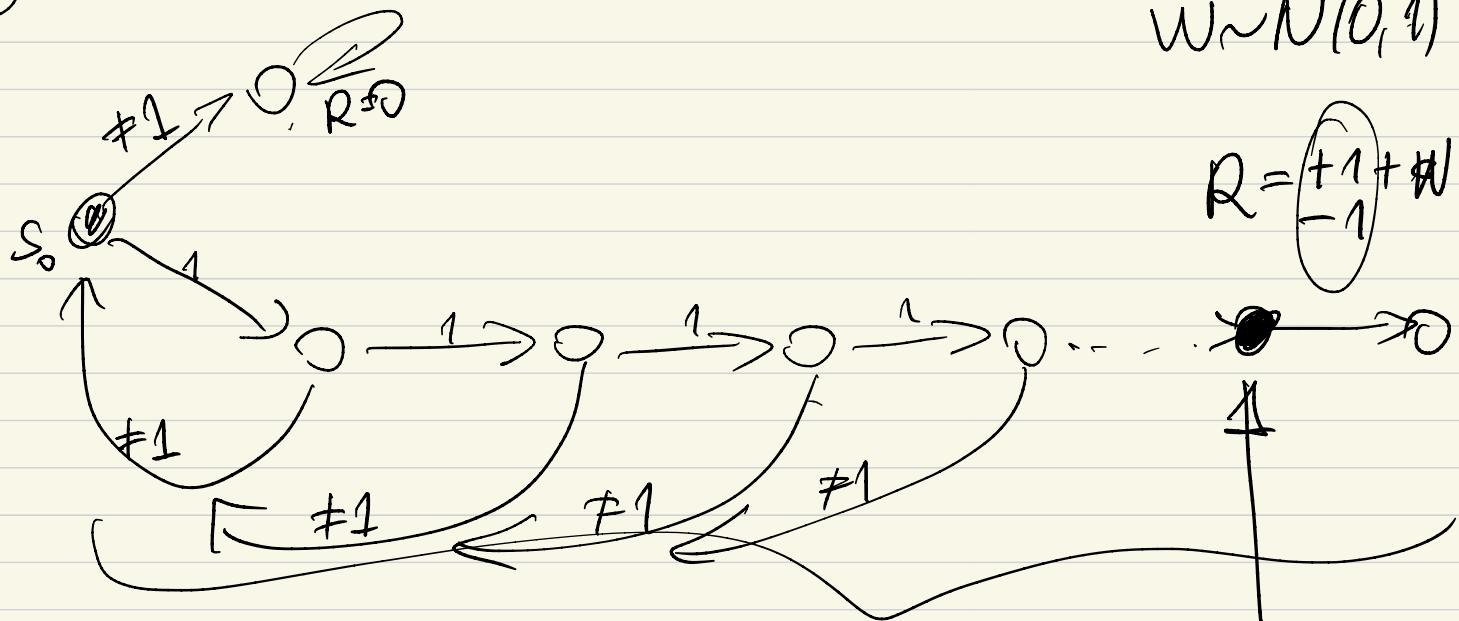
$$= 1 - \tau - \frac{1}{3} \frac{(1-\tau)}{\tau} = (1-\tau) \cdot \frac{2}{3}$$

$$(1 - P(\tau)) P(\tau) = \frac{1}{3} \frac{1-\tau}{\tau} \cdot \frac{4\tau-1}{3\tau} \underset{\substack{\text{=} \\ \text{---}}}{=} \approx \frac{1}{3}(1-\tau) \underset{\sim}{\overset{\tau \rightarrow 1}{\longrightarrow}} 1$$

$$\frac{\frac{\gamma}{(1-\gamma p)^2}}{\underbrace{\sqrt{\frac{p(1-p) \log(\dots)}{2n}}}_{\approx}} \approx \frac{1}{\underbrace{(1-\gamma)^2 \left(\frac{2}{3}\right)^2}_{\frac{\frac{1}{2}(1-\gamma) \log(\dots)}{2n}}} \leq \varepsilon$$

$$n \geq H^3 / \varepsilon^2 \log(\dots) \quad H = \frac{1}{1-\gamma}$$

3r



$$\pi_b(1|s)$$

$$H = H_{\sigma, \varepsilon} \quad P_0^2 \left(\frac{1}{A}\right)^H$$

$$\leq \pi_b(a|s) \quad \text{fa}$$

$$\boxed{\pi_b(1|s) \leq \frac{1}{A}}$$

$$\boxed{n \geq P_0^{-1} = AH}$$