

Contextual Bandits LUCB

- ϵ -Greedy

- Square CB (inverse prob. gap weighting)

Proved given Lemma $\forall \delta > 0$

$$\forall \hat{f}, f^*: [K] \rightarrow [0, 1]$$

$$\pi \sim P(\pi) = \frac{1}{\lambda + \alpha (\hat{A}(\hat{\pi}) - \hat{f}(\pi))}$$

$$\mathbb{E}_{\rho} [f^*(\pi^*) - f^*(\pi)]$$

$$\leq \frac{K}{\gamma} + \gamma \mathbb{E} [(f(\pi) - f^*(\pi))^2]$$

Pf: $\mathbb{E}_{\rho} [f^*(\pi^*) - f^*(\pi)]$

$$= \mathbb{E}_{\rho} [f^*(\pi^*) - \hat{f}(\pi^*)$$

$$+ \hat{f}(\pi^*) - \hat{f}(\pi)$$

+ ...

]

$$\leq \frac{K}{\sigma} + \sigma \mathbb{E}_{\mathcal{P}} \left[(f(\pi) - f^*(\pi))^2 \right]$$

$\mathbb{E} \text{ Regret} \leq \frac{KT}{\sigma} + \sigma \sum_{t=1}^T \mathbb{E}_{\mathcal{P}} \left[(F_t(\pi) - f^*(\pi))^2 \right]$
 (over T rounds)

Regret sq

$$\sigma = \sqrt{\frac{KT}{\text{Regret sq}}}$$

$$= 2 \sqrt{KT \text{Regret sq}}$$

Today: Lin UCB

High-level: ① Use linear regression ($\hat{\theta}_t$)
to estimate the reward
of each arm

② Use uncertainty estimates $\hat{\theta}_t$
to get UCB on each arm

③ "Optimistic" pick
arm w/ best UCB.

Algorithm: $\theta^* \in \mathbb{R}^d$, $\beta > 0$ (Assume $\phi(\cdot, \cdot)$ known)

For $t = 0$ to $T-1$.

$$\|\theta^*\|_2 \leq 1$$

- $\hat{\theta}_t = \underset{\|\theta\|_2 \leq 1}{\operatorname{argmin}} \sum_{s=1}^t (r_s - \underbrace{\langle \theta, \phi(x_s, \pi_s) \rangle}_{\phi_s})^2$

$\mathbb{E}[r_t | x_t, \pi_t] = \langle \theta^*, \phi(x_t, \pi_t) \rangle$

- Define

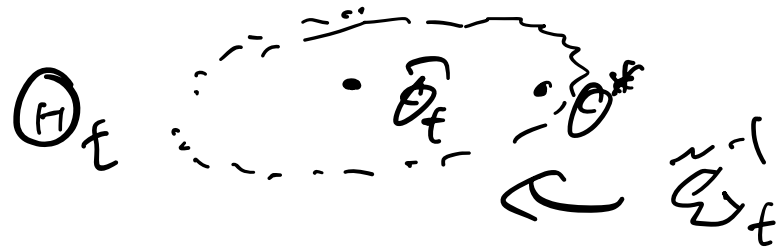
$$\tilde{\Sigma}_t = \sum_{s=1}^t \phi_s \phi_s^T + \mathbf{I}$$

$$\text{(so } \langle u, \tilde{\Sigma}_t^{-1} u \rangle = \sum_{s=1}^t \langle \phi_s, u \rangle^2 + \|u\|_2^2)$$

"how much we know about direction u "

- Observe context x_{t+1} .

- Select



$$\pi_{t+1} \in \operatorname{argmax}_{\pi \in [K]} \left[\begin{array}{l} \max_{\theta} \langle \theta, \phi(x_{t+1}, \pi) \rangle \\ \theta \text{ s.t.} \end{array} \right]$$

uncertainty set

$$\Theta_t$$

$$\langle \theta - \hat{\theta}_t, \tilde{z}_t \rangle$$

$$\leq (16\beta + 4)^2$$

- Play π_{t+1} , get reward r_{t+1} .

Possible issues:

- Θ_t too small, doesn't include θ^*

- Θ_t too big.

Thm. $\mathbb{E} \text{Regret} = \mathbb{E} \sum_{t=1}^T (f^*(x_t, \pi_t^*) - f^*(x_t, \pi_t))$

$$\lesssim d \sqrt{T} \log T$$

where $\beta = \Theta(d \log T)$.

PF Has two parts.

Part I: Show (with high probability)

x^* is always in the uncertainty set.

pf:

$$\pi_{t+1}^* = \operatorname{argmax}_{\pi \in \mathcal{K}} \langle \theta^*, \phi(x_{t+1}, \pi) \rangle$$

$$\mathbb{E} \text{ Regret} = \sum_{t=0}^{T-1} \left(\langle \theta^*, \phi(x_{t+1}, \pi_{t+1}^*) \rangle - \langle \theta^*, \phi(x_{t+1}, \pi_{t+1}) \rangle \right)$$

Step 1 (FW) $\theta^* \in \Theta_t$

$$\leq \sum_{t=0}^{T-1} \left(\max_{\theta \in \Theta_t} \langle \theta, \phi(x_{t+1}, \pi_{t+1}^*) \rangle - \langle \theta^*, \phi(x_{t+1}, \pi_{t+1}) \rangle \right)$$

$$\leq \sum_{t=0}^{T-1} \left(\max_{\theta \in \Theta_t} \langle \theta, \phi(x_{t+1}, \pi_{t+1}) \rangle - \langle \theta^*, \phi(x_{t+1}, \pi_{t+1}) \rangle + \max_{\theta \in \Theta_t} \langle \theta - \theta^*, \phi(x_{t+1}, \pi_{t+1}) \rangle \right)$$

$$\max_{\Theta} \left(\Theta - \Theta^*, \phi(x_t, \pi_t) \right)$$

$$= \max_{\Theta} \left(\sum_t \Theta^{1/2} (\Theta - \Theta^*), \sum_t \phi(x_t, \pi_t)^{-1/2} \right)$$

Lemma (Elliptical Potential Lemma)

Suppose $\phi_1, \dots, \phi_T \in \mathbb{R}^D$ $\|\phi_i\| \leq 1$

$$\sum_t \phi_t = \sum_{s=1}^T \phi_s \phi_s^\top + I$$

Then

$$\sum_{t=1}^T \phi_t \sum_{t=1}^{t-1} \phi_t \leq 2D \log T.$$

"how surprising
 ϕ_t is"

Why inverse?

$$Y = X\Theta^* + \text{noise}$$

$$X^{-1}Y = \Theta^* + X^{-1} \text{noise}$$

(classic OLS)

Completion of proof of Thm, given Lemma.

$$\mathbb{E} \text{Regret} \leq \mathbb{E} \sum_{t=1}^T \max_{\theta \in \Theta_{t-1}} \langle \theta - \theta^*, \phi(x_t, \pi_t) \rangle$$

$$\leq \mathbb{E} \sum_{t=1}^T \max_{\theta \in \Theta_{t-1}} \left\langle \sum_{t=1}^{t-1} \frac{1}{2} (\theta - \theta^*), \sum_{t=1}^{t-1} \phi(x_t, \pi_t) \right\rangle$$

$$\leq \mathbb{E} \sum_{t=1}^T \sqrt{\langle \theta - \theta^*, \sum_{t=1}^{t-1} (\theta - \theta^*) \rangle} \sqrt{\langle \phi(x_t, \pi_t), \sum_{t=1}^{t-1} \phi(x_t, \pi_t) \rangle}$$

$$\leq \sqrt{16\beta^2 + 4} \sqrt{T} \sqrt{\sum_{t=1}^T \langle \phi(x_t, \pi_t), \sum_{t=1}^{t-1} \phi(x_t, \pi_t) \rangle}$$

Pf of Elliptical potential lemma.

Key idea: Compute $\det \tilde{\Sigma}_{t+1}$
in terms of $\det \tilde{\Sigma}_t$.

Factorise remaining matrices

$$\textcircled{1} \det(M) = \prod_{i=1}^d \lambda_i(M)$$

$$\textcircled{2} \lambda_i(I + uu^T)$$

$$= 1 + u^T u$$

$$\det(\tilde{\Sigma}_{t+1}) = \det\left(\tilde{\Sigma}_t + \phi_{t+1} \phi_{t+1}^T\right)$$

$$= \det\left(\tilde{\Sigma}_t^{1/2} \left(I + \tilde{\Sigma}_t^{-1/2} \phi_{t+1} \phi_{t+1}^T \tilde{\Sigma}_t^{-1/2} \right)\right)$$

$$= \det(\tilde{\Sigma}_t) \det\left(I + \tilde{\Sigma}_t^{-1/2} \phi_{t+1} \phi_{t+1}^T \tilde{\Sigma}_t^{-1/2} \right)$$

$$= \det(\tilde{\Sigma}_t) \left(1 + \phi_{t+1}^T \tilde{\Sigma}_t^{-1} \phi_{t+1} \right)$$

$$\log \det \tilde{\Sigma}_{t+1} - \log \det \tilde{\Sigma}_t$$

$$= \log \left(I + \Phi_{t+1}^T \tilde{\Sigma}_t^{-1} \Phi_{t+1} \right)$$

$$\geq \frac{\Phi_{t+1}^T \tilde{\Sigma}_t^{-1} \Phi_{t+1}}{2}$$

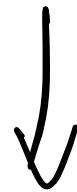
$$\log \det \tilde{\Sigma}_T - \log \det \tilde{\Sigma}_0 \geq \sum_{t=0}^{T-1} \frac{\Phi_{t+1}^T \tilde{\Sigma}_t^{-1} \Phi_{t+1}}{2}$$

$\log \det I = 0$

$$\sum_{i=1}^d \log \lambda_i(\tilde{\Sigma}_T) \leq d \log(T+1)$$

Fact:
 $\log(1+x) \geq \frac{x}{2}$ for $x \in [0, 1]$

'total amount of surprise'



Picture of what happens during Lin VCR



$$\textcircled{H} \Leftrightarrow \Sigma^{-1}, \quad \textcircled{\Theta} = 0$$

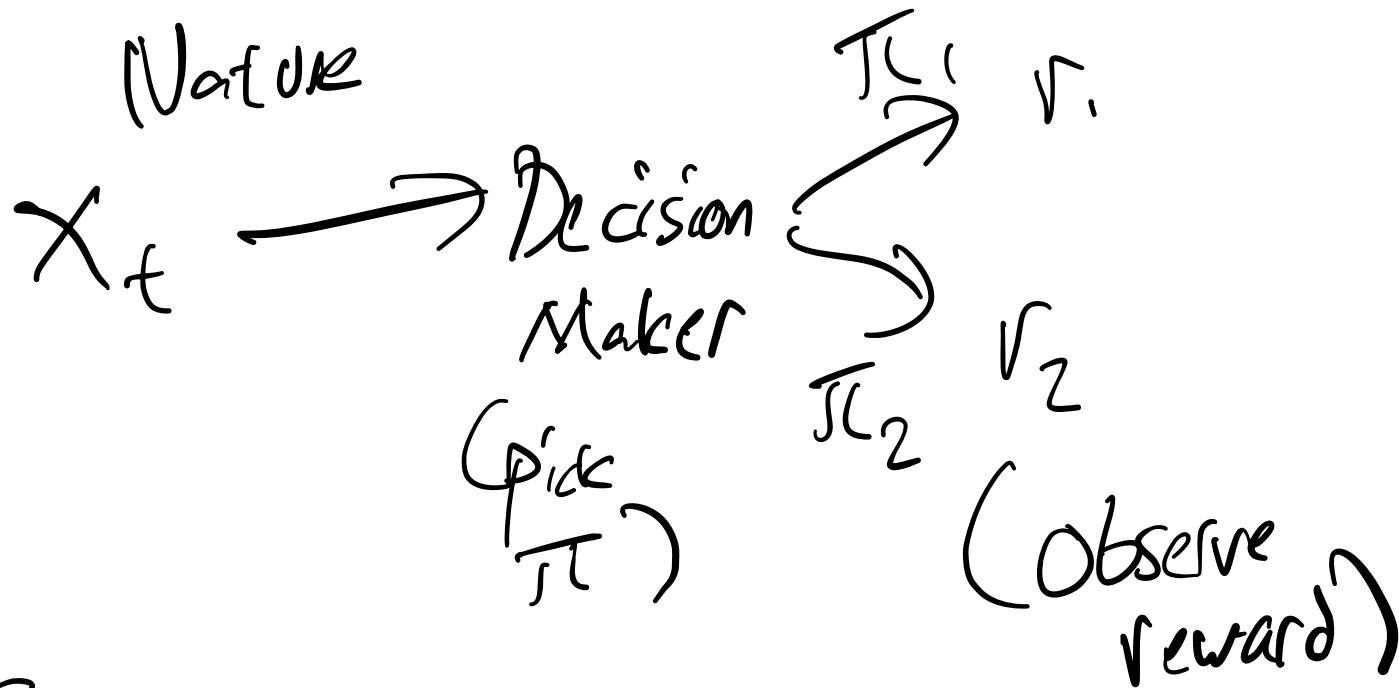
$$\textcircled{\Theta}_1 \rightarrow \phi_1 = e_1 \rightarrow \Sigma_1^{-1} = \Sigma_0^{-1} + \phi_1 \phi_1^T$$

$$\rightarrow \phi_2 = e_2$$

$$\text{Vol ellipse} \Leftrightarrow \det \Sigma_t^{-1} \\ = \frac{1}{\det \Sigma_t}$$

Markov Decision Processes (Intro to RL)

Sofar: bandit, contextual bandit



Regret \longleftrightarrow for the same seq of X_t ,
how good vs. best decision in hindsight

What is an MDP?

$$S = \{ \text{state space} \}$$

$$A = \{ \text{action space} \}$$

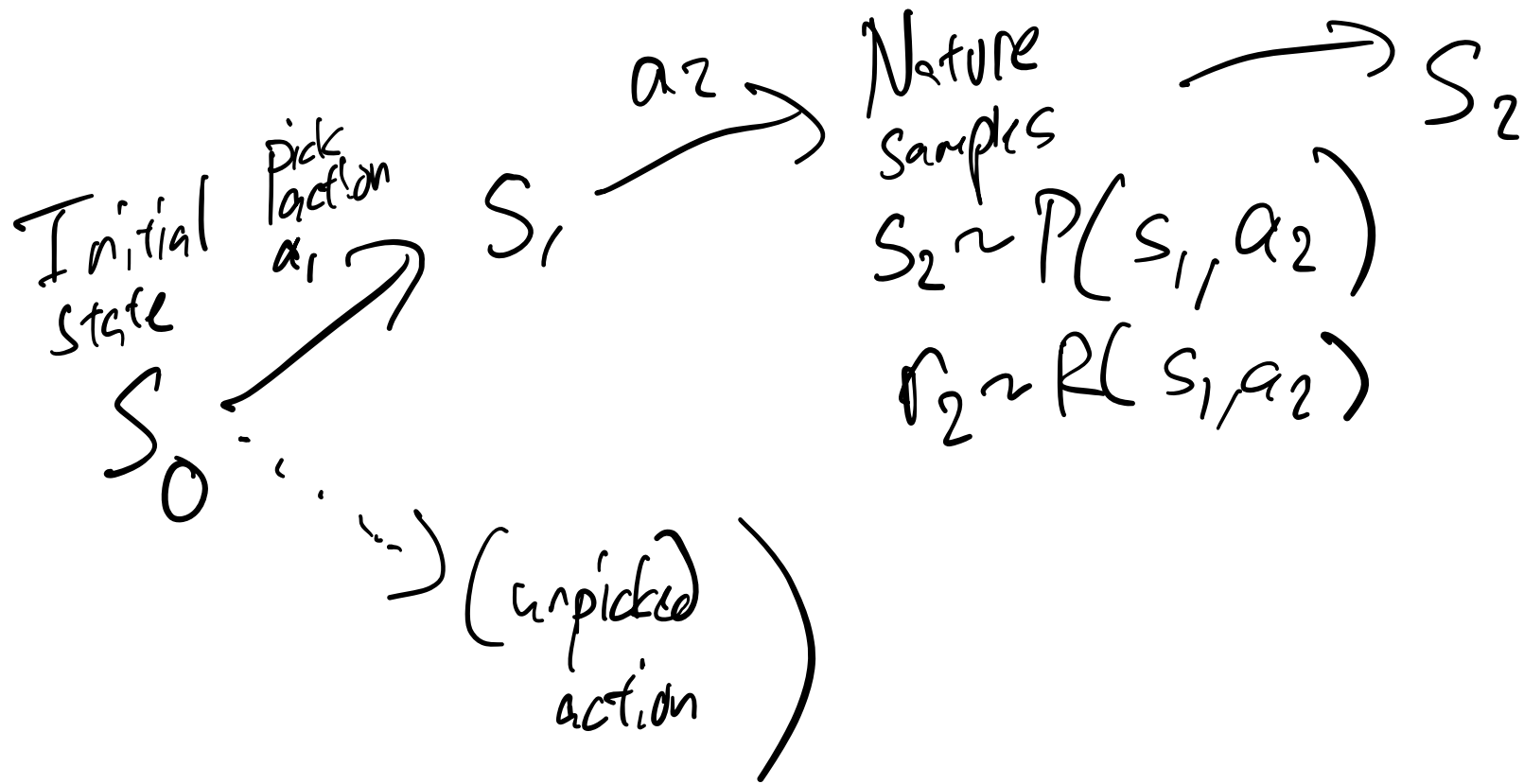
$$\Delta(S) = \{ \text{probdists over } S \}$$

"Transition kernel"

$$P: S \times A \rightarrow \Delta(S)$$

"Reward distribution"

$$R: S \times A \rightarrow \Delta(\mathbb{R})$$



Goal: Optimize total reward.