

Stochastic Linear Bandits

• Stochastic contextual bandits:

- similar (mirror) of prev. lecture

- key difference,

$$X_t = \underbrace{r(C_t, A_t)}_{\substack{\text{reward} \\ \text{function}}} + \underbrace{\eta_t}_{\text{sub-Gaussian noise}}$$

sub-Gaussian noise

$$r: \mathcal{C} \times [k] \rightarrow \mathbb{R}$$

- w/o going into ~~sign~~ σ -fields, idea is that

$$E[X_t | H_{t-1}] = r(C_t, A_t)$$

- if $r(\cdot, \cdot)$ was known, then,

$$A_t^* \in \operatorname{argmax}_{a \in [k]} r(C_t, a)$$

and so, the expected regret is

$$R_n = E \left[\sum_{t=1}^n \max_a r(C_t, a) - \sum_{t=1}^n X_t \right]$$

Note: since no assumption is made on how rewards are chosen, it is possible to pick sub-optimal arms, one way to circumvent is to ensure that A_1, \dots, A_{t+1} don't significantly alter C_t, \dots, C_n 's.

• feature map-based

- assume learner has access to $\psi: C \times [k] \rightarrow \mathbb{R}^d$,
and then, for some unknown $\theta^* \in \mathbb{R}^d$,

$$r(c, a) = \langle \theta^*, \psi(c, a) \rangle \quad \forall (c, a)$$

- ψ is referred to as a feature map
[O/P of NN before final layer]

$$|r(c, a) - r(c', a')| \leq \|\theta^*\| \|\psi(c, a) - \psi(c', a')\|_2$$

[Hölder]

so, assumptions on $\theta^* \equiv$ "smoothness" of $r(\cdot, \cdot)$

• Stochastic Linear Bandits

- the "action" is not as critical as the feature vector of an action, and thus,

- let $A_t \subset \mathbb{R}^d$ be decision action set, $A_t \in A_t$

$$X_t = \langle \theta^*, A_t \rangle + \eta_t$$

then,

$$\hat{R}_n = \sum_{t=1}^n \max_{a \in A_t} \langle \theta^*, a \rangle$$

$$R_n = E \left[\sum_{t=1}^n \max_{a \in A_t} \langle \theta^*, a \rangle - X_t \right]$$

"Decision set"

- if $A_t = \{e_1, \dots, e_d\}$, then stochastic bandit
- if $A_t = \{\psi(E_t, i) \mid i \in [K]\}$, then contextual linear bandit.

How to solve?

- Generalization of UCB: pick $C_t \subset \mathbb{R}^d$ that "contains θ_* w.h.p."

- assume, $\mathbb{E}_t \theta_* \in C_t$, then for any $a \in \mathbb{R}^d$, let

$$UCB_t(a) = \max_{\theta \in C_t} \langle \theta, a \rangle$$

be an upper bound on mean payoff $\langle \theta^*, a \rangle$;
then select

$$A_t = \operatorname{argmax}_{a \in \mathcal{A}_t} UCB_t(a)$$

- but how to construct C_t ??
- need an estimator for θ_* [analogous to μ_i for Stoch]

- regularized least-squares estimator.

$$\hat{\theta}_t = \operatorname{argmin}_{\theta \in \mathbb{R}^d} \left(\sum_{s=1}^t (x_s - \langle \theta, A_s \rangle)^2 + \lambda \|\theta\|^2 \right)$$

$$\hat{\theta}_t = V_t^{-1} \sum_{s=1}^t A_s x_s ; \quad V_0 = \lambda I ; \quad V_t = V_0 + \sum_{s=1}^t A_s A_s^T$$

now, just pick

$$C_t \subseteq E_t = \left\{ \theta \in \mathbb{R}^d \mid \|\theta - \hat{\theta}_{t+1}\|_{V_{t+1}}^2 \leq \beta_t \right\}$$

Remarks : A. $\|x\|_D^2 = x^T D x$, $D \succeq 0$

B $\beta_1 \geq 1, \dots, \beta_1 \leq \beta_2 \leq \dots \leq \beta_n$

C: $\lambda_i(V_t)$ are "increasing" \Rightarrow $\text{vol}(E_t)$ is shrinking

Main Result :

Assumption : w' follo win hold

- $1 \leq \beta_1 \leq \beta_2 \leq \dots \leq \beta_n$

- $\max_{t \in [n]} \sup_{a, b \in A_t} \langle \theta_{t+1}, a-b \rangle \leq 1$

- $\|a\|_2 \leq L \quad \forall a \in \bigcup_{t=1}^n A_t$

- $\exists \delta \in (0, 1)$ s.t w'p $1-\delta$, $\forall t \in [n]$
 $\theta_{t+1} \in C_t$

Thm : w'p $1-\delta$, lin VCB satisfies

$$\hat{R}_n \leq \sqrt{8n \beta_n \log \left(\frac{\det V_n}{\det V_0} \right)} \leq \sqrt{8dn \beta_n \log \left(\frac{d\lambda + nL^2}{d\lambda} \right)}$$

$$R_n \leq Cd\sqrt{n} \log(nL)$$

$$\sqrt{\beta_n} = \sqrt{\lambda} \|\theta_{t+1}\|_2 + \sqrt{2 \log(1/\delta) + d \log \left(\frac{d\lambda + nL^2}{d\lambda} \right)}$$