

6.2 Linear Bandits - Handling Large Action Spaces

Let D be a compact set of decision. On each round, we must choose $x_t \in D$, each of which results in a reward r_t .

Pseudo code : The Linear UCB algorithm

Input : λ, β_t

1: for $t=0, 1, \dots$ do

2: Execute $x_t = \operatorname{argmax}_{x \in D} \max_{\mu \in \text{BALL}_t} \mu \cdot x$

and observe r_t

3: Update BALL_{t+1}

4: End for.

Note:

① We assume $E[r_t | x_t = x]$ is a fixed linear function,
i.e. $E[r_t | x_t = x] = \mu^* \cdot x \in [-1, 1]$ for all $x \in D$

② We assume $E[r_t] \in [-1, 1]$. Then the **noise sequence**
 $\eta_t = r_t - \mu^* \cdot x_t$ is a martingale difference seq.

③ If x_0, \dots, x_{T-1} are the decisions made in the game,
the cumulative regret is defined by

$$R_T = T \mu^* \cdot x^* - \sum_{t=0}^{T-1} \mu^* \cdot x_t$$

where $x^* \in D$ is an optimal decision for μ^* ,

i.e. $x^* \in \operatorname{argmax}_{x \in D} \mu^* \cdot x$

x^* exists since D is compact

④ Our goal is to keep R_T as small as possible.

1. The LinUCB algorithm.

A. pseudo code : we define an upper bound vector BALL_t

at episode t , we define an uncertain region $BALL_t$.

The center of $BALL_t$ is $\hat{\mu}_t$, which is the solution of:

$$\begin{aligned}\hat{\mu}_t &= \arg \max_{\mu} \sum_{T=0}^{t-1} \|\mu \cdot x_T - r_T\|_2^2 + \lambda \|\mu\|_2^2 \\ &= \Sigma_t^{-1} \sum_{T=0}^{t-1} r_T x_T\end{aligned}$$

where λ is the parameter and Σ_t satisfies:

$$\Sigma_t = \lambda I + \sum_{T=0}^{t-1} x_T x_T', \quad \Sigma_0 = \lambda I$$

Then $BALL_t$ is defined as

$$BALL_t = \{ \mu \mid (\hat{\mu}_t - \mu)' \Sigma_t (\hat{\mu}_t - \mu) \leq \beta_t \}$$

2. Upper and Lower bounds

Assume $x \in \mathbb{R}^d$

Theorem 6.3. Suppose that the expected costs are bounded by 1, i.e. $|\mu^* \cdot x| \leq 1$ for all $x \in D$; $\|\mu^*\| \leq W$, $\|x\| \leq B$ for all $x \in D$, and that η_t is σ^2 sub-Gaussian. Set

$$\lambda = \sigma^2 / W^2, \quad \beta_t := \sigma^2 \left(2 + 4d \log \left(1 + \frac{tB^2W^2}{d} \right) + 8 \log(4/\delta) \right)$$

we have that with prob. $\geq 1 - \delta$, for all $T \geq 0$

$$R_T \leq c \sigma \sqrt{T} \left(d \log \left(1 + \frac{TB^2W}{d\sigma^2} \right) + \log(4/\delta) \right),$$

where c is an absolute constant. i.e. $R_T \sim O(d\sqrt{T})$

Following shows R_T of Th.6.3 is best.

Theorem 6.4 [lower bound] There exist a distribution over linear bandit problems (i.e. $\Delta(\mu)$) with rewards bounded by 1, in martingale, and $\sigma^2 \leq 1$, s.t. for every algorithm, we have for $n \geq \max\{256, d^2/16\}$,

$$E_{\mu} E[R_T] \geq \frac{1}{2500} d \sqrt{T}$$

w.r.t. randomness in the problem & algorithm.

We will eliminate the dependencies in Th 6.3.

Let Σ_D denote the D-optimal design matrix from :

Th 3.2: Suppose $\mathcal{X} \subset \mathbb{R}^d$ is a compact set. There exists a distribution p on \mathcal{X} s.t.

① p is supported on at most $d(d+1)/2$ points

② Define $\Sigma = E_{x \sim p} [xx^T]$, we have
$$\|x\|_{\Sigma^{-1}}^2 \leq d$$

and p is referred to as the D-optimal design.

Coordinate transformation:

$$\tilde{x} = \Sigma_D^{-1/2} x, \quad \tilde{\mu}^* = \Sigma_D^{1/2} \mu^*$$

$$\text{s.t. } \tilde{x} \cdot \tilde{\mu}^* = x' \Sigma_D^{-1/2} \cdot \Sigma_D^{1/2} \mu^* = x \cdot \mu^*$$

we hold the expected reward function.

$$\text{We have } \underbrace{\|\tilde{x}\|^2}_{\uparrow B} = \|x\|_{\Sigma_D^{-1}}^2 \leq d \quad (\text{by Th. 3.2})$$

$$\text{and } \underbrace{\|\tilde{\mu}^*\|}_{\uparrow W} = \|\mu^*\|_{\Sigma_D} = \sqrt{(\mu^*)' \Sigma_D \mu^*} = \sqrt{E_{x \sim p} [(\mu^* \cdot x)^2]} \stackrel{(\text{assume } \|r\| \leq 1)}{\leq} 1$$

Following shows that we could remove the dependencies on B and W from previous theorem, due to $B \leq \sqrt{d}$ and $W \leq 1$.

Corollary 6.5 Suppose that the expected rewards are bounded in martingale by 1, $\|\tilde{\mu}^* \cdot \tilde{x}\| \leq 1$ for all $x \in \mathcal{D}$; and that η_t is σ^2 sub-Gaussian. Suppose linUCB is implemented in the \tilde{x} coordinate system, with following settings:

$$\lambda = \sigma^2, \quad \beta_{ti} = \sigma^2 (2 + 4d \log(1+t) + 8 \log(4/s))$$

with prob. $\geq 1 - \delta$, for all $T \geq 0$,

$$R_T \leq c \sigma \sqrt{T} (d \log(1 + \frac{T}{\sigma^2}) + \log(4/s))$$