

1) Setting: Linear Contextual Bandits

- contexts $x \in \mathcal{X} \subseteq \mathbb{R}^d$
drawn from distribution $\mathbb{D} \in \Delta(\mathcal{X})$ $x_t \stackrel{iid}{\sim} \mathbb{D}$
- action "arms" $a \in \mathcal{A} = \{1, \dots, K\}$
- rewards $r_t = r(x_t, a_t)$
 $\mathbb{E}[r(x, a)] = \mu_a(x) = \Theta_a^T x$ linear function
 $\Theta_a \in \mathbb{R}^d$ (unknown)
- horizon T

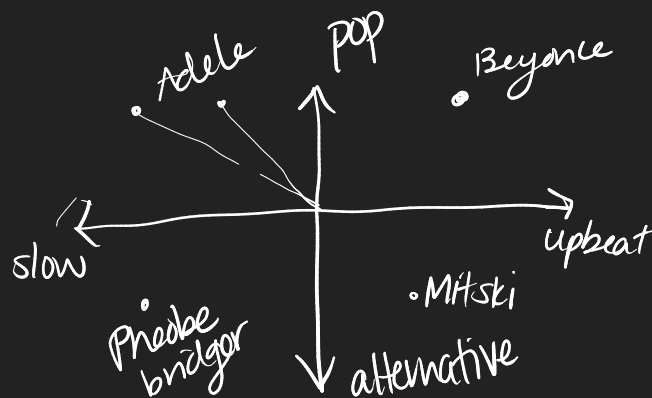
Goal: find a policy $a_t = \pi(x_t)$, achieve low regret

$$R(T) = \sum_{t=1}^T \mathbb{E} \left[\underbrace{\max_a \Theta_a^T x_t}_{\mu_*(x_t)} - \underbrace{\Theta_{a_t}^T x_t}_{\pi_*(x_t)} \right]$$

Example: music recommendation

arm = artists

context = user's affinity towards traits



Last lecture: explore-then-commit based on supervised learning

$$a_t = \operatorname{argmax}_a \hat{\mu}_a(x_t)$$

$$\hat{\mu}_a = \operatorname{argmin}_{\mu \in \mathcal{M}} \sum_{i=1}^N (\mu(x_i^a) - r_i^a)^2$$

↑ data during exploration ↑ collected exploration

Linear Regression

since $\eta_a(x) = \theta_a^T x$, $M = \{ \eta(x) = \theta^T x \mid \theta \in \mathbb{R}^d \}$

$$\hat{\theta}_a = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N (\theta^T x_i^a - r_i^a)^2$$

Lemma: as long as $(x_i)_{i=1}^N$ span \mathbb{R}^d

$$\hat{\theta} = \underbrace{\left(\sum_{i=1}^N x_i x_i^T \right)}_A^{-1} \underbrace{\sum_{i=1}^N x_i r_i}_b = A^{-1} b$$

Proof: $\nabla_{\theta} \left[\sum_{i=1}^N (\theta^T x_i - r_i)^2 \right] = 2 \sum_{i=1}^N (\theta^T x_i - r_i) x_i = 0$

$$\underbrace{\left[\sum_{i=1}^N x_i x_i^T \right]}_A \theta = \underbrace{\sum_{i=1}^N x_i r_i}_b$$

□

A is related to empirical covariance of x

$$\Sigma = \mathbb{E}[x x^T] \quad \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N x_i x_i^T \quad A = N \hat{\Sigma}$$

3) LinUCB Algorithm

We keep track of:

$$A_t^a = \sum_{k=1}^t x_k x_k^T \mathbb{1}_{\{a_k = a\}} \quad b_t^a = \sum_{k=1}^t x_k r_k \mathbb{1}_{\{a_k = a\}}$$

$$\hat{\theta}_t^a = (A_t^a)^{-1} b_t^a$$

Alg: LinUCB

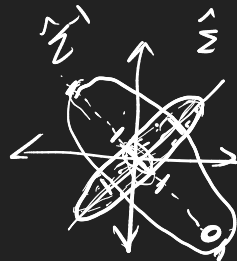
Initializing 0 mean & infinite CI

For $t=1, \dots, T$:

observe x_t

$$a_t = \operatorname{argmax}_a \left[\hat{\Theta}_t^{aT} x_t + \alpha \sqrt{x_t^T (A_t^a)^{-1} x_t} \right]$$

observe r_t , update $\hat{\Theta}^{a_t}$, b^{a_t} , A^{a_t}



Intuition

$$\hat{\Theta}^T x + \alpha \sqrt{x^T A^{-1} x}$$

large if $\hat{\Theta}$ is aligned w/ x

large when past data is not aligned with x

$$A = N \hat{\Sigma}$$

$$x^T A^{-1} x = \frac{1}{N} \underbrace{x^T \hat{\Sigma}^{-1} x}$$

amount of data

large when x is not aligned with prev. data

Statistical Explanation

claim: with high probability

$$\Theta_a^T x \leq \hat{\Theta}_a^T x + \alpha \sqrt{x^T A_a^{-1} x}$$

where α depends on probability and variance of rewards

Lemma (Chebychev's Inequality)

For a random variable u with $\mathbb{E}[u] = 0$,

$$|u| \leq \beta \sqrt{\mathbb{E}(u^2)}$$

with probability $1 - \frac{1}{\beta^2}$

Proof of Claim:

We will show, using Chebychev's

$$\underbrace{|\hat{\Theta}_a^T x - \Theta_a^T x|}_{=u} \leq \alpha \sqrt{x^T A_a^{-1} x}$$

from computing $\mathbb{E}[u]$

$$1) \mathbb{E}[\hat{\theta}^T x - \theta^T x] \stackrel{?}{=} 0$$

$$\text{let } w_i = r_i - \underbrace{\mathbb{E}r_i}_{\theta^T x_i}$$

$$\hat{\theta} = \left(\sum_{i=1}^N x_i x_i^T \right)^{-1} \sum_{i=1}^N x_i (\theta^T x_i + w_i)$$

$$\hat{\theta} = \left(\sum_{i=1}^N x_i x_i^T \right)^{-1} \left[\left(\sum_{i=1}^N x_i x_i^T \right) \theta + \sum_{i=1}^N x_i w_i \right]$$

$$= \theta + \left(\sum_{i=1}^N x_i x_i^T \right)^{-1} \sum_{i=1}^N x_i w_i$$

$$\mathbb{E}[\hat{\theta}^T x - \theta^T x] = \mathbb{E}[(\hat{\theta} - \theta)^T x] = \mathbb{E} \left[\left(\sum_{i=1}^N x_i x_i^T \right)^{-1} \sum_{i=1}^N x_i w_i \right]^T x$$

$$\mathbb{E} \left[A^{-1} \sum_{i=1}^N x_i w_i \right]$$

$$= A^{-1} \sum_{i=1}^N x_i \mathbb{E}[w_i]$$

$$\mathbb{E}[w_i] = 0$$

$$= 0$$

✓