

# Lecture 20: Linear Contextual Bandits

## 1) Setting

Our simplified MDP setting consists of:

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

Goal: find a policy  $a_t = \pi(x_t)$  that achieves low regret.

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

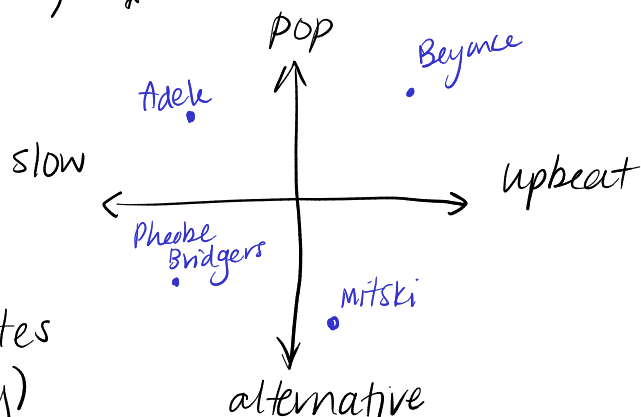
Example: music recommendation

arms  $a$  are artists

$\Theta_a \in \mathbb{R}^d$  represents attributes

$x \in \mathbb{R}^d$  represents a users

affinity towards the attributes  
(observed from listening history)



Last lecture we considered an explore-then-commit algorithm for general function approximation/supervised learning

$$a_t = \arg \max_a \hat{y}_a(x_t) \quad \text{where}$$

$$\hat{y}_a = \arg \min_{\mu \in M} \sum_{i=1}^N (\mu(x_i^a) - r_i^a)^2$$

↑  
data collected during exploration phase.

## Linear Regression

If we know that  $\mu_a(x) = \theta_a^T x$  we can instantiate the general supervised learning framework with

$$M = \{ \mu(x) = \theta^T x \mid \theta \in \mathbb{R}^d \}$$

In this case the learning problem is equivalent to

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

We will sometimes drop the  $a$  subscript in these notes.

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

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

Proof:

$$\nabla_{\theta} \sum_{i=1}^N (\theta^T x_i - r_i)^2 = 2 \sum_{i=1}^N x_i (x_i^T \theta - r_i)$$

setting the gradient equal to zero,

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

The matrix on the left hand side is invertible if  $(x_i)_{i=1}^N$  span  $\mathbb{R}^d$

(why? Let  $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \in \mathbb{R}^{N \times d}$ . Then if  $x_i$  span  $\mathbb{R}^d$ ,  $X$  has full row rank,  $\text{rank}(X) = d$

$\sum_{i=1}^N x_i x_i^T = X^T X \in \mathbb{R}^{d \times d}$  is full rank because  $\text{rank}(X^T X) = \text{rank}(X) = d$ . Therefore it is invertible. )

□

The matrix  $A$  is related to the empirical covariance

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

We can relate  $A = N \hat{\Sigma}$ .

2) Interactive Demo - dyputer Notebook

### 3) LinUCB Algorithm

Recall that last lecture we wanted to estimate conditional errors  $\mathbb{E}[(\hat{y}_a(x) - y_a(x))^2 | x]$ . Using the structure of the linear regression problem, we can do this.

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

Initialize 0 mean & infinite confidence intervals

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

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

update  $\hat{\Theta}_t^{a_t}, b_t^{a_t}, A_t^{a_t}$

#### Geometric Intuition:

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

large if  $x$  and  $\hat{\Theta}$  are aligned

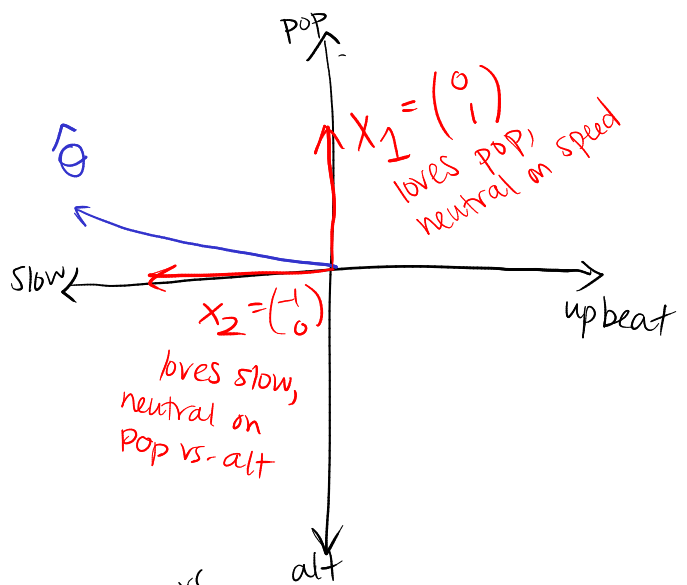
large if  $x$  is not aligned with much historical data

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

amount of data

alignment w/ data

ex-  $\hat{\Theta}^T x_2 > \hat{\Theta}^T x_1$   
 ↓ ↓  
 large alignment on speed      less alignment on pop



ex- if previous data is

$(x_1, -x_1, x_1, x_1, -x_1, -x_2)$

then  $A^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1/5 \end{bmatrix}$   
 ↗ ↘  
 less sure on speed      pretty sure about pop

most users so far are on neutral speed

$$\hat{\Sigma} = \begin{bmatrix} 1/6 & 0 \\ 0 & 5/6 \end{bmatrix}$$

$x_1^T A^{-1} x_1 < x_2^T A^{-1} x_2$   
 ↗ ↘  
 more sure      less sure

## Statistical Explanation:

Claim: With high probability (over noisy rewards)

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

where  $\alpha$  depends on probability & variance of rewards

Lemma: (Chebychev's Inequality)

for a random variable  $u$  with  $\mathbb{E}(u) = 0$ ,

$$|u| \leq \beta \sqrt{\mathbb{E}(u^2)} \text{ with probability } 1 - 1/\beta^2$$

Proof: we will use chebychev's to show that w.h.p

$$|\underbrace{\hat{\theta}_a^T x - \theta_a^T x}_u| \leq \alpha \sqrt{\underbrace{x^T A^{-1} x}_{\mathbb{E}u^2}}$$

1) compute expectation. Define  $w_i = r_i - \mathbb{E}[r_i]$  so  $r_i = \theta_a^T x_i + w_i$ .

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

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

$$\theta - \theta_a = A^{-1} \sum_{i=1}^N x_i w_i$$

$$\text{therefore } \mathbb{E}((\theta - \theta_a)^T x) = A^{-1} \sum_{i=1}^N x_i \mathbb{E}[w_i] = 0$$

2) compute variance

$$\begin{aligned} \mathbb{E}_{\theta} \left[ ((\theta - \theta_a)^T x)^2 \right] &= \mathbb{E}_{w_i} \left[ x^T A^{-1} \sum_{i=1}^N x_i w_i \cdot \sum_{i=1}^N x_i^T w_i A^{-1} x \right] \\ &= x^T A^{-1} \mathbb{E} \left[ \sum_{i=1}^N \sum_{j=1}^N x_i x_j^T \underbrace{w_i w_j}_{\text{}} \right] A^{-1} x \end{aligned}$$

The noise in rewards is iid so the expectation is 0 if  $i \neq j$ . Define  $\sigma^2$  as variance of rewards.

$$= x^T A^{-1} \sum_{i=1}^N x_i x_i^T \sigma^2 A^{-1} x$$

$\underbrace{\quad}_{=A}$

$$= \sigma^2 x^T A^{-1} x$$

Therefore, using Chebychev's, we have that w.p.  $1 - 1/\beta^2$ ,

$$|\theta_a^T x - \hat{\theta}_a^T x| \leq \underbrace{\beta \sigma}_{\alpha} \sqrt{x^T A^{-1} x}$$

Thus the upper bound of this confidence interval is  $\hat{\theta}^T x + \alpha \sqrt{x^T A^{-1} x}$

□