

# Contextual Bandits

Subrahmanya Swamy Peruru

# Contextual Bandits – Multiple states

- News article Recommendation systems



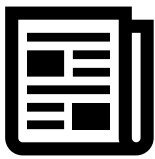
- Articles – arms
- Like / Dislike – Reward
- User – State

Different users have different preferences for articles

# Multi-arm Bandits – One state

Arms /  
Articles


*a*



Expected

$$\mu(\underline{a}) = \underline{0.2}$$

*b*


$$\mu(\underline{b}) = \underline{0.9}$$

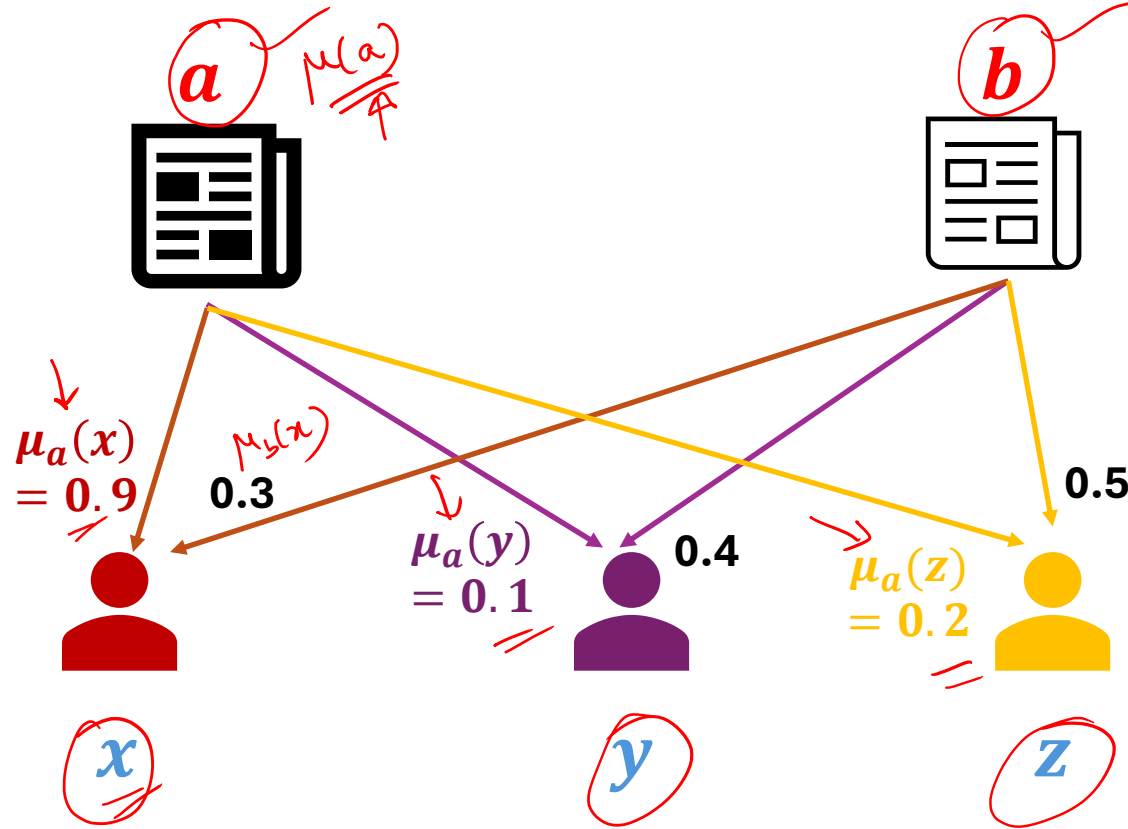
- Each arm has only one expected reward associated with it

# Contextual Bandits – Multiple States

ETC  
 $\epsilon$ -greedy  
UCB

Arms /  
Articles

Users /  
States

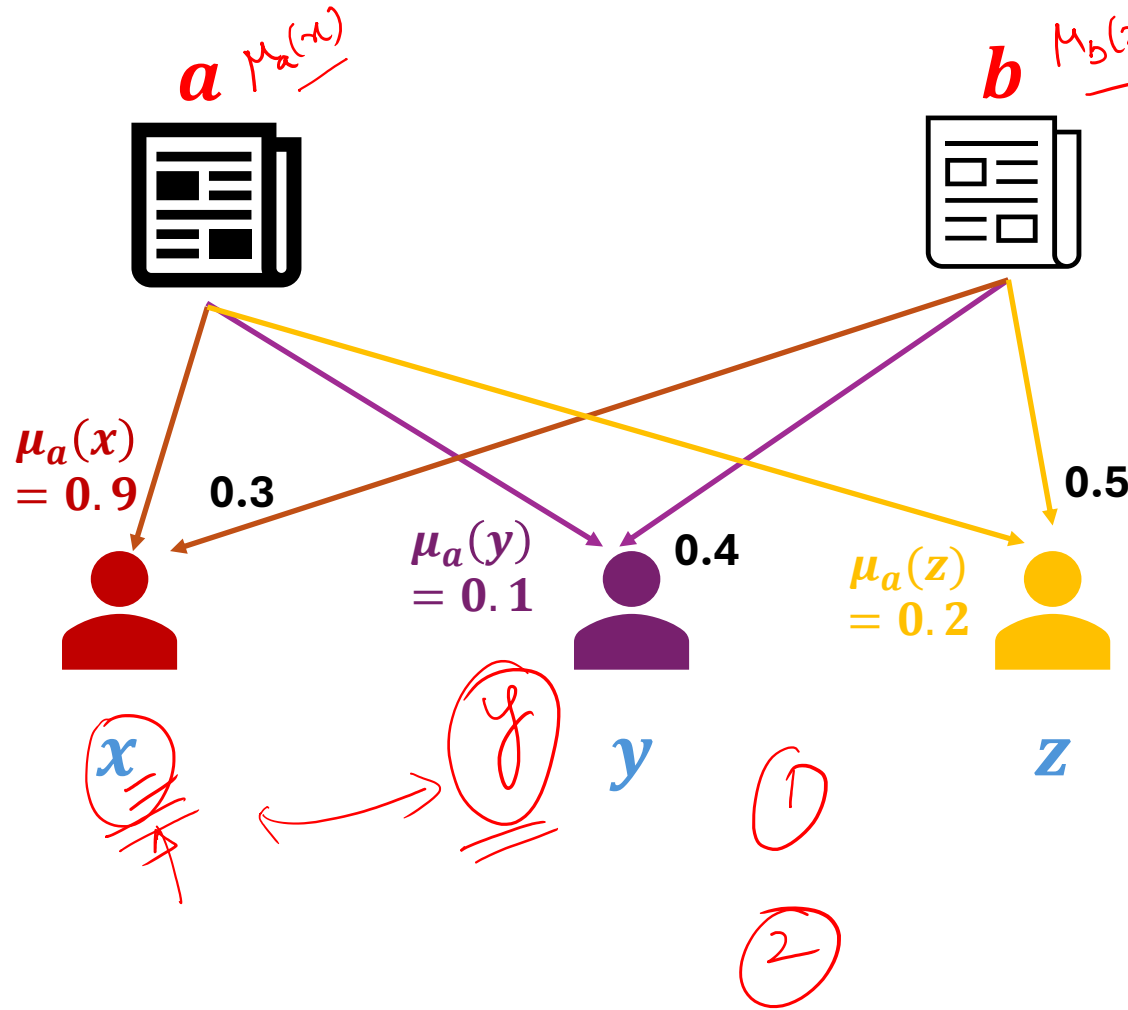


- Expected reward of an arm changes with user  $\mu_a(x)$
- How to deal with it?**

# Contextual Bandits – Multiple States

Arms /  
Articles

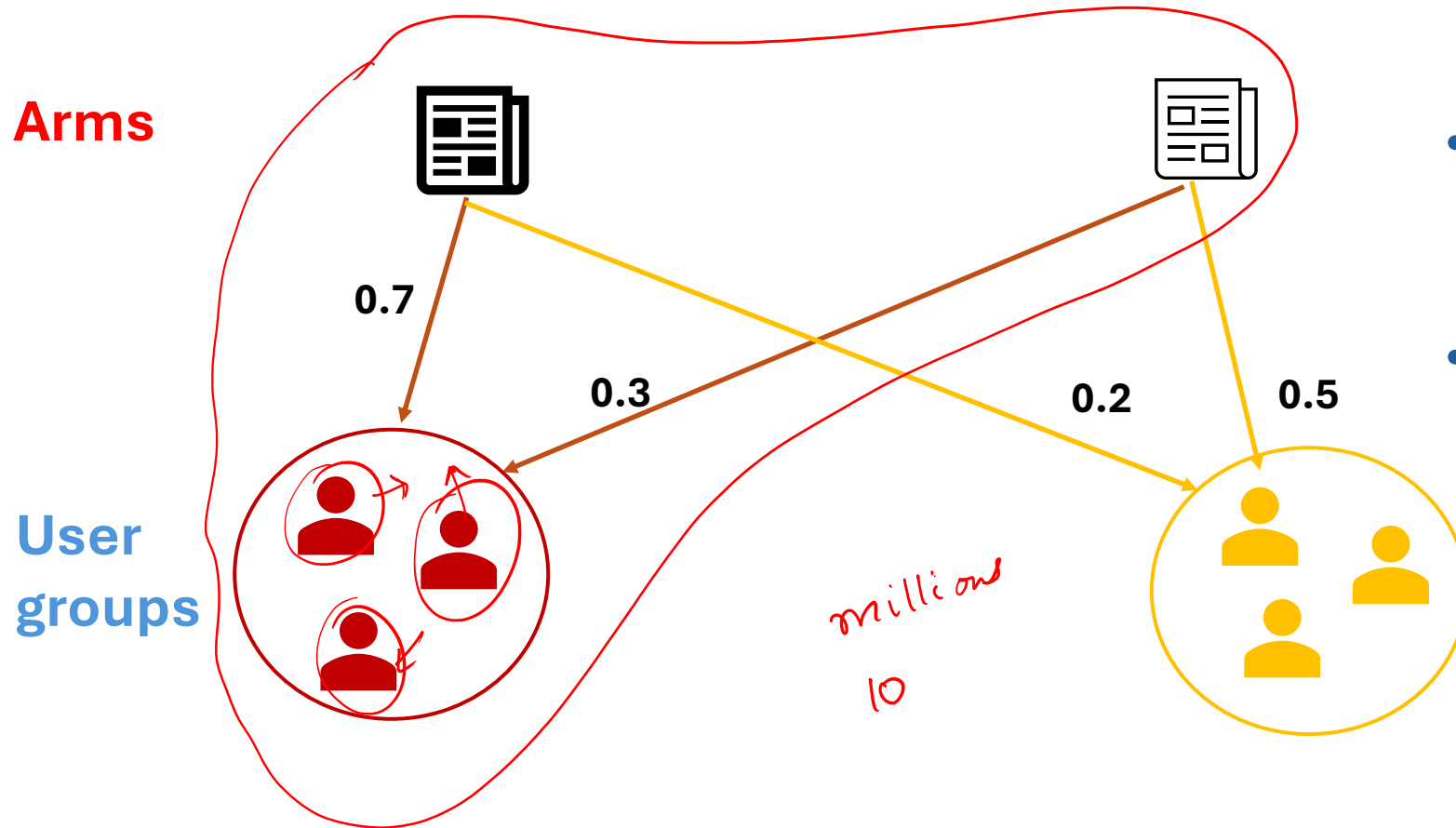
Users /  
States



unsup  
sup

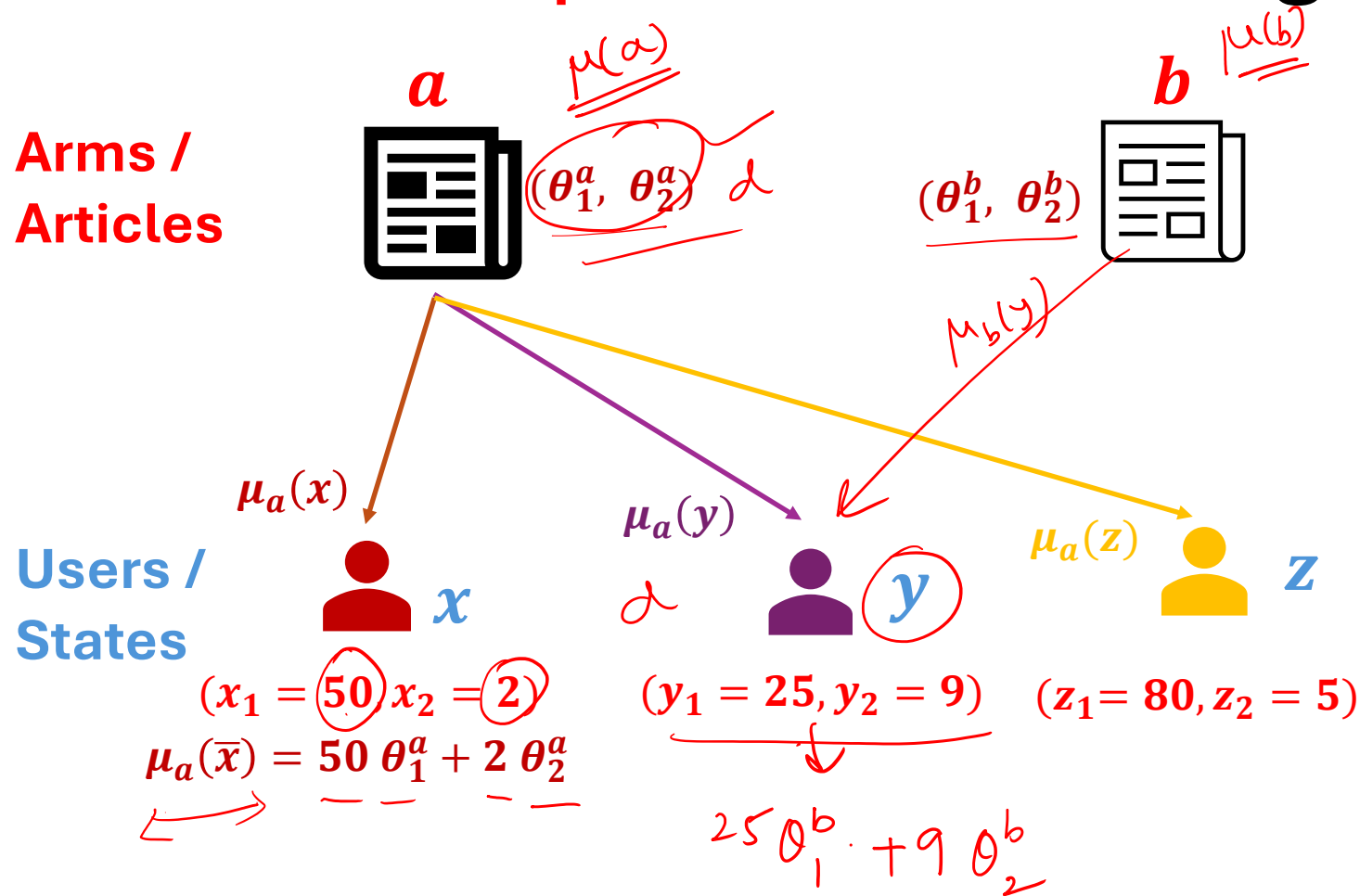
- Expected reward of an arm changes with user  $\mu_a(x)$
- Treat each user as a separate bandit problem
- Practically infeasible with millions of users!

# Bandits + Unsupervised Learning



- Form user groups by clustering similar users together
- Solve a separate bandit problem for each cluster

# Bandits + Supervised Learning



ETC

- User (state) represented by features such as age, income  
 $\bar{x} = (x_1, x_2)$
- Model the expected reward for user  $\bar{x}$  for pulling arm  $a$  as  
 $\mu_a(\bar{x}) = \theta_1^a x_1 + \theta_2^a x_2$

# Contextual (Linear) Bandits

- Users (state) represented by features such as age, gender
  - State feature vector:  $\underline{\bar{x}} = (x_1, x_2)^T$  Eg: (50, 1), (25, 0), (80, 1)
- Each article (arm) has a different expected reward associated with each user (state)
  - The expected reward of an arm is characterized by unknown  $\bar{\theta}^a = (\theta_1^a, \theta_2^a)^T$
  - State-specific expected reward:  $\underline{\mu_a(\bar{x})} = \theta_1^a x_1 + \theta_2^a x_2$  (Linear Bandits)
- The reward for playing arm  $a_t$  under state  $\underline{\bar{x}_t}$  is  $R_t = \underline{\mu_{a_t}(\bar{x}_t)} + \epsilon_t$ , where  $\epsilon_t$  is independent mean-zero noise.
- If "T arbitrary users" visit our website sequentially, How to maximize the total reward  $\sum_{t=1}^T R_t$ ?

$$\begin{array}{l} \mu(a) = 2 \\ \downarrow \\ R_1 = \begin{array}{l} 2 \cdot 1 \\ 2 \cdot 7 \\ 1 \cdot 8 \end{array} \left. \vphantom{\begin{array}{l} 2 \cdot 1 \\ 2 \cdot 7 \\ 1 \cdot 8 \end{array}} \right\} \frac{2 \cdot 1 + 2 \cdot 7 + 1 \cdot 8}{3} \end{array}$$



# Solution Approach

- Explore each arm N times and
- Estimate the mean parameters based on sample rewards

## One-state Multi-arm Bandits

$\mu(a) = 2$  **Unknown parameter**

Sample reward:  $R_t = \underline{\mu(a)} + noise$

Rewards from **3 rounds of exploration**

$$\begin{aligned} R_1 &= \mu(a) \approx 2.7 \\ R_2 &= \mu(a) \approx 1.6 \\ R_3 &= \mu(a) \approx 2.1 \end{aligned}$$

Best estimate:  $\hat{\mu}(a) = \arg \min_x \underbrace{(R_1 - x)^2 + (R_2 - x)^2 + (R_3 - x)^2}$

# Linear Bandits

Unknown parameter:  $(\theta_1^a, \theta_2^a) = \theta_a$

ETC  
↓ a → 3 times

Sample reward:  $R_t = \mu_a(x) + \text{noise}$   
 $= \theta_1^a x_1 + \theta_2^a x_2 + \text{noise}$

## Exploration:

3 users with features:  $(50, 1)$ ,  $(25, 4)$ ,  $(80, 7)$

Observed rewards:  $0.7, 0.4, 0.5$

$$\begin{aligned} R_1 &= 0.7 \approx \theta_a^1 \underline{50} + \theta_a^2 \underline{1} \\ R_2 &= 0.4 \approx \theta_a^1 \underline{25} + \theta_a^2 \underline{4} \\ R_3 &= 0.5 \approx \theta_a^1 \underline{80} + \theta_a^2 \underline{7} \end{aligned}$$

✓

$D_a$  ✓

$\theta_a$

$\theta_a = \begin{bmatrix} \theta_a^1 \\ \theta_a^2 \end{bmatrix}$

$\begin{bmatrix} 50 & 1 \\ 25 & 4 \\ 80 & 7 \end{bmatrix} \begin{bmatrix} \theta_a^1 \\ \theta_a^2 \end{bmatrix} = \begin{bmatrix} 0.7 \\ 0.4 \\ 0.5 \end{bmatrix}$

$b_a$  ✓

Best estimate:  $\hat{\theta}^a = \arg \min_{(\theta_1^a, \theta_2^a)} \left\{ \begin{aligned} & (0.7 - \theta_a^1 \underline{50} + \theta_a^2 \underline{1})^2 + \\ & (0.4 - \theta_a^1 \underline{25} + \theta_a^2 \underline{4})^2 + \\ & (0.5 - \theta_a^1 \underline{80} + \theta_a^2 \underline{7})^2 \end{aligned} \right\}$

↑ ↑

$\min_{\theta_1^a, \theta_2^a} f(\theta_1^a, \theta_2^a)$

Linear Regression

# ETC algorithm for Contextual Bandits

- Explore each arm  $N$  times
- Based on the history  $\{(\vec{x}_t, a_t, R_t)\}_{t=1}^{NK}$ , estimate the parameters for all  $a \in \mathcal{A}$  using Ridge regression.

- $\{\hat{\theta}^a = (D_a^T D_a + I_d)^{-1} D_a^T b_a\}$
- $D_a$  is  $N \times d$  context matrix whose rows represent user feature vectors
- $b_a$  is  $N \times 1$  reward vector with rewards obtained during  $N$  exploration rounds

$D_a = \begin{bmatrix} x_1^1 & x_2^1 \\ x_1^2 & x_2^2 \\ x_1^3 & x_2^3 \end{bmatrix}$  feature dimension  $d = 2$ , arm  $a$  played  $N = 3$  times

$NK$

- At any round  $t > \underline{NK}$ , play the arm  $a_t = \arg \max_a x_t^T \hat{\theta}^a$

$a \rightarrow (\hat{\theta}_a^1, \hat{\theta}_a^2) \checkmark$   
 $b \rightarrow (\hat{\theta}_b^1, \hat{\theta}_b^2) \checkmark$

$x_t \rightarrow (x_t^1, x_t^2)$   
 so,

$$\hat{\mu}(a) = \frac{\bar{\mu}(a)}{\bar{\mu}(b)} \quad \frac{N}{K} \quad \max \{ \bar{\mu}(a), \bar{\mu}(b) \}$$

$\epsilon$ -greedy

$0.2$   
 $\downarrow$   
 $x_t^1 \hat{\theta}_b^1 + x_t^2 \hat{\theta}_b^2$

$0.9$   
 $\downarrow$   
 $(x_t^1 \hat{\theta}_a^1 + x_t^2 \hat{\theta}_a^2)$

$$\text{UCB score} = \underbrace{\bar{\mu}_t(a)}_{\text{exploit}} + \underbrace{\sqrt{\frac{2 \log T}{n_t(a)}}}_{\text{exploration}}$$

$$\hat{\mu}_a(x_t) +$$

$(1,2)$   
 $(2,3)$

$(1,2)$   
 $(1,2)$   
 $(1,2)$

$\epsilon$ -greedy

$$\left. \begin{aligned} \hat{\mu}_t(a) &\approx \bar{\mu}(a) \\ \hat{\mu}_t(b) &\approx \bar{\mu}(b) \end{aligned} \right\}$$

$t:$   
 prob:  $\epsilon \rightarrow \text{heads} \rightarrow \text{Explore}$   
 $1-\epsilon \rightarrow \text{Tail} \rightarrow \text{Exploit}$

$$\begin{aligned} &(\hat{\theta}_a^1(t), \hat{\theta}_a^2(t)) \\ &(\hat{\theta}_b^1(t), \hat{\theta}_b^2(t)) \end{aligned}$$

$t$   
 prob  $\epsilon \rightarrow \text{heads} \rightarrow \text{random } a, b$   
 $1-\epsilon \rightarrow \text{exploit}$   
 $x_t \rightarrow \hat{\mu}_a(x_t) = \frac{\hat{\theta}_a^1(t)x_t^1 + \hat{\theta}_a^2(t)x_t^2}{\hat{\mu}_b(x_t)}$

# LinUCB (UCB for Linear Bandits)

- At time  $t$  based on the history  $\{(\bar{x}_s, \overset{\downarrow}{a_s}, \overset{\uparrow}{R_s})\}_{s=1}^{t-1}$ , **estimate** the parameters for all  $a \in \mathcal{A}$  using **Ridge regression**.
  - $\widehat{\theta}^a = (D_a^T D_a + I_d)^{-1} D_a^T b_a$
  - If arm  $a$  is played  $m$  times till round  $t$ ,
  - $D_a$  is  $m \times d$  context matrix whose rows correspond to the respective user feature vectors
  - $b_a$  is  $m \times 1$  reward vector corresponding to the rewards obtained during those rounds

$\widehat{D}_a = \begin{bmatrix} x_1^1 & x_2^1 \\ x_1^2 & x_2^2 \\ x_1^3 & x_2^3 \end{bmatrix}$  feature dimension  $d = 2$ , arm  $a$  played  $m = 3$  times

- Pick the arm  $a_t = \arg \max_a \underbrace{x_t^T \widehat{\theta}^a}_{\text{Exploit}} + \underbrace{\sqrt{x_t^T (D_a^T D_a + I_d)^{-1} x_t}}_{\text{Explore}}$