

## 1) Setting: Contextual Bandits

Simplified RL setting: simplified version of state: context  
context is memoryless, drawn iid at each timestep

$\mathcal{X}$ : set of contexts  $x$

$\mathcal{A} = \{1, \dots, k\}$  a set of discrete actions

$\mathbb{D} \in \Delta(\mathcal{X})$ : context distribution  $x_t \stackrel{iid}{\sim} \mathbb{D}$

$r: \mathcal{X} \times \mathcal{A} \rightarrow \Delta(\mathbb{R})$  noisy reward  $r_t \sim r(x_t, a_t)$   
depends on context & action  
 $\mathbb{E}[r(x, a)] = \mu_a(x)$

$T$ : time horizon

Actions should depend on the context,  
 $\pi: \mathcal{X} \rightarrow \mathcal{A}$  (or  $\pi$  (als) stochastic)

Optimal Policy:  $\pi^*(x) = \operatorname{argmax}_a \mu_a(x)$

Minimize Expected Regret:

$$R(T) = \sum_{t=1}^T \mathbb{E}_{x_t \sim \mathbb{D}} [\mu^*(x_t) - \mu_{a_t}(x_t)]$$

$$\mu^*(x_t) = \max_a \mu_a(x_t)$$

## 2) Tabular Setting

Suppose  $M$  contexts

IDEA: run a separate MAB algorithm for each context

Alg: Explore-then-Commit w/ Context

For  $t=1, 2, \dots, T$

Observe  $x_t$

1) if  $\exists$  arm pulled less than  $N$  times — explore  
for context  $x_t$ , pull it

2) otherwise,  $a_t = \operatorname{argmax}_a \hat{\mu}_a(x_t)$  — exploit

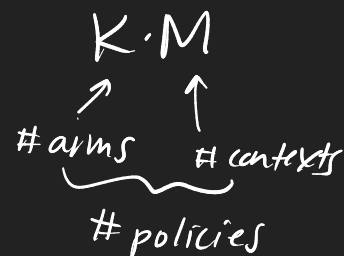
# Alg: VCB w/ contexts

For  $t=1, 2, \dots, T$ :

$$\text{pull } a_t = \underset{a}{\operatorname{argmax}} \hat{\mu}_t^a(x_t) + \sqrt{\frac{\log(TKM/s)}{N_t^a(x_t)}}$$

↑  
context-dependent  
mean & count

K·M policies: Regret bounds will be similar to last 2 lectures with K replaced w/



## 3) Function Approximation

We may never see the same context twice!

ex: user 1: { F, 22, CS } =  $x_1$

user 2: { M, 21, econ } =  $x_2$

user 3: { F, 20, econ } =  $x_3$

Instead of estimating  $\hat{\mu}_a(x)$  with counting we can use function approximation

$$\hat{\mu}_a(x) = \underset{\mu \in \mathcal{M}}{\operatorname{argmin}} \sum_{k=1}^t (\mu(x_k) - r_k)^2 \mathbb{1}\{a_t = a\}$$

↑  
function class

How to get CI on  $\hat{\mu}_a(x)$ ?  
Error bounds for supervised learning

Lemma: for  $x_i \sim \mathcal{D}$ ,  $\mathbb{E}[y_i] = f_*(x_i)$   $f_* \in \tilde{\mathcal{F}}$

$$\hat{f} = \operatorname{argmin}_{f \in \tilde{\mathcal{F}}} \sum_{i=1}^N (\hat{f}(x_i) - y_i)^2$$

Then with high probability,

$$\mathbb{E}_{x \sim \mathcal{D}} [|\hat{f}(x) - f_*(x)|] \lesssim \sqrt{\frac{C_{\tilde{\mathcal{F}}}}{N}} \leftarrow \text{complexity of } \tilde{\mathcal{F}}$$

Algorithm: Explore-then-Commit w/ SL

- 1) pull each arm  $N$  times, record  $\{ \{ (x_i^a, r_i^a) \}_{i=1}^N \}_{a=1}^K$   
( $t=1, \dots, NK$ )  
Estimate  $\hat{\mu}_a(x) = \operatorname{argmin}_{\mu} \sum_{i=1}^N (\mu(x_i^a) - r_i^a)^2$
- 2)  $t = NK+1, \dots, T$ : pull  $a_t = \operatorname{argmax}_a \hat{\mu}_{a_t}(x_t)$

Regret Analysis:

$$R(T) = \underbrace{R_1}_{\leq NK} + \underbrace{R_2}_T = \sum_{NK+1}^T \mathbb{E}_{x_t \sim \mathcal{D}} [\mu_{a^*}(x_t) - \mu_{a_t}(x_t)]$$

$$\mathbb{E} [\mu_{a^*}(x_t) - \mu_{a_t}(x_t)] = \mathbb{E} \left[ (\mu_{a^*}(x_t) - \hat{\mu}_{a^*}(x_t)) + \hat{\mu}_{a^*}(x_t) - \hat{\mu}_{a_t}(x_t) + \hat{\mu}_{a_t}(x_t) - \mu_{a_t}(x_t) \right] \leq 0$$

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

$$\leq \mathbb{E}_{x_t \sim \mathcal{D}} [|\mu_{a^*}(x_t) - \hat{\mu}_{a^*}(x_t)|] + \mathbb{E}_{x_t \sim \mathcal{D}} [|\hat{\mu}_{a_t} - \mu_{a_t}(x_t)|]$$

$$\lesssim 2 \sqrt{\frac{C_M}{N}}$$

$$R(T) \lesssim \overset{p_1 \downarrow}{N} K + \overset{p_2 \downarrow}{2T} \sqrt{\frac{C_M}{N}} \quad N = \left( \frac{T}{2K} \sqrt{C_M} \right)^{2/3}$$

$$R(T) \lesssim T^{2/3} (K C_M)^{1/3}$$

UCB algorithm?

Naive:  $\hat{\mu}_a^t(x) + \sqrt{\frac{C_M}{N_a^t}}$

We want confidence intervals based on conditional expected error

$$\mathbb{E}[|y(x) - \hat{y}(x)| | x]$$

Next lecture: Lin UCB algorithm

$$y_a(x) = \theta_a^T X$$

General contextual setting

$$y_a(x) = \theta_a^T \phi(x, a)$$