

Lecture 19: Contextual Bandits

1) Motivation - slides

2) Formal Setting

Simplified RL setting with simplified version of state: context. Unlike states, contexts are memoryless. They are drawn from a fixed distribution independent of previous context & actions.

\mathcal{X} : a set of contexts x

$A = \{1, \dots, k\}$ 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 : integer time horizon

Notice that the context distribution \mathbb{D} in some sense subsumes the transition probabilities P and the initial distribution μ_0

The actions should depend on the observed context, therefore policy $\pi: \mathcal{X} \rightarrow \mathcal{A}$ (or stochastic $\pi: \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$)

The optimal policy maximizes cumulative reward at each step:

$$\pi^*(x) = \max_a \mu_a(x)$$

Q: why is it sufficient to consider each timestep independently?
How is this different from full MDP?

Goal: Minimize expected Regret (in terms of cumulative reward)

$$R(T) = \sum_{t=1}^T \mathbb{E}_{x_t \sim \mathbb{D}} \left[\underbrace{\max_a \mu_a(x_t)}_{\mu_{a^*}(x_t)} - \mu_{a_t}(x_t) \right]$$

3) Naive (Tabular) Approach

Suppose there are M discrete possible contexts.

IDEA: run a separate MAB for each different context.

Instead of computing mean-per-arm $\hat{\mu}_a$, compute mean-per-arm-and-context

$$\hat{\mu}_a(x) = \frac{\sum_{k=1}^t r_k \mathbb{1}\{a_k = a\} \mathbb{1}\{x_k = x\}}{\sum_{k=1}^t \mathbb{1}\{a_k = a\} \mathbb{1}\{x_k = x\}}$$

number of times a pulled in context x

Algorithm: Explore - then - Commit with context:

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

observe x_t

If $\exists a$ s.t. # times a pulled in context $x_t \leq N$ } exploration

$a_t = a$

else:

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

} exploitation } context dependent

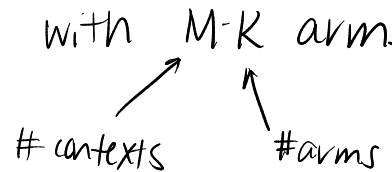
Algorithm: UCB with contexts:

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

pull $a_t = \operatorname{argmax}_a \hat{\mu}_t^a(x_t) + \sqrt{\frac{\log(CKTM/S)}{N_t^a(x_t)}}$

Keep track of confidence intervals for arm-context pairs.

This is similar to a classic MAB but with $M \cdot K$ arms.
 We can show similar regret bounds as previous lectures where $K \leftarrow MK$.



In some sense we are searching over all possible $M \cdot K$ policies (vs. K actions).

4) Function Approximation

In reality, contexts include many pieces of information (e.g. demographic information, recent browsing behaviour, etc) and the number of discrete contexts may be very large! We may never see the exact same context twice!

However, it is also likely that correlations exist between similar contexts

ex: user 1: {gender: F, age: 22, major: CS} = x_1
user 2: {gender: M, age: 21, major: econ} = x_2
user 3: {gender: F, age: 21, major: econ} = x_3

Information about user 1 & user 2 should help us predict for user 3.

Instead of estimating $\hat{y}_a(x)$ by counting, we can use function approximation:

$$\hat{y}_a(x) = \underset{\text{function class}}{\operatorname{argmin}}_{y \in M} \frac{\sum_{k=1}^t (y(x_k) - r_k)^2 \mathbb{1}\{a_k = a\}}{\sum_{k=1}^t \mathbb{1}\{a_k = a\}}$$

Q: how to get confidence intervals on $\hat{y}_a(x)$?

A: supervised learning guarantees:

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

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

We have that, with high probability,

$$\mathbb{E} [|\hat{f}(x) - f_*(x)|] \lesssim \sqrt{\frac{c_{\mathcal{F}}}{n}}$$

Algorithm: Explore-then-commit w/ fn. approx

1) pull each arm N times and record $\{ \{ x_i^a, r_i^a \}_{i=1}^N \}_{a=1}^K$
 Estimate $\hat{\mu}_a(x) = \underset{\mu \in M}{\operatorname{argmin}} \sum_{i=1}^N (\mu(x_i^a) - r_i^a)^2$

2) For $t = NK + 1, \dots, T$ pull $a_t = \underset{a}{\operatorname{argmax}} \hat{\mu}_a(x_t)$

Regret analysis:

$$R(T) \leq \underbrace{NK}_{\text{from explore}} + \sum_{NK+1}^T \mathbb{E}_{x_t \sim D} [\mu_{a^*}(x_t) - \mu_{a_t}(x_t)]$$

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

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

Then $R(T) \lesssim NK + 2T \sqrt{\frac{CM}{N}}$ very similar to non-contextual!

$$\text{Set } N = \left(\frac{T}{2K} \sqrt{CM} \right)^{2/3}$$

$$\text{so that } R(T) \lesssim T^{2/3} (KC_M)^{1/3}$$

What about UCB type algorithm? Good Confidence intervals require knowing conditional expected error

$$\mathbb{E}[\mu_a(x) - \hat{\mu}_a(x) \mid x]$$

Next lecture: Linear contextual bandits & LinUCB

$$\mu_a(x) = \theta_a^\top x$$