

Stochastic Contextual Bandits

Source: Bandit Algorithms, L&S, Chap 19.

Setting:

At time $t = 1, 2, \dots, n$:

1. Player observes context $C_t \in \mathcal{C}$.

2. Player plays action $A_t \in \mathcal{A}_t$

3. Nature chooses reward $X_t = r(C_t, A_t) + \eta_t$

$r: \mathcal{C} \times \mathcal{A}_t \rightarrow \mathbb{R}$, a real valued reward function.

$\eta_t \sim 1\text{-subGaussian noise.}$

$I_t = (C_1, A_1, X_1, C_2, A_2, X_2, \dots, C_t, A_t)$

$E[e^{\lambda \eta_t} | I_t] \leq e^{\lambda^2/2}$ a.s. (i.e., conditionally
1-subG noise).

3@: r : a known function.

4. $R_n(\pi, v) = E \left[\sum_{t=1}^n \max_{a \in \mathcal{A}_t} r(C_t, a) - \sum_{t=1}^n X_t \right]$

(Note:

Back to stochastic setting: At each time, a different action can potentially be chosen).

$r(\cdot, \cdot)$ through a Feature Map:

$$r(c, a) = \langle \theta^*, \psi(c, a) \rangle, \quad \theta^* \in \mathbb{R}^d,$$

\downarrow
feature map $\in \mathbb{R}^d$

example: c = visitor to a website; one-hot encoding.
 a = ad to place; one-hot encoding.

Leads to Linear Bandit; directly work in feature space.

$$\text{Model: } X_t = \langle \theta^*, A_t \rangle + \eta_t$$

unknown vector $\in \mathbb{R}^d$.

} action by player $\in A_t$

$$\mathbb{E}[\eta_t | I_t] \leq e^{\gamma^2/2}, \quad \text{where}$$

$$I_t = \{A_1, A_1, X_1, \dots, A_{t-1}, A_{t-1}, X_{t-1}, A_t, A_t\}$$

$$R_n = E \left[\sum_{t=1}^n \max_{a \in A_t} \langle \theta^*, a \rangle - \sum_{t=1}^n x_t \right]$$

(a) $A_t = \{e_1, e_2, \dots, e_d\}$ unit vectors in \mathbb{R}^d .

Then,

$$\langle \theta^*, e_i \rangle = \theta_i^*, \text{ i.e.,}$$

$$x_t = \theta_i^* + \eta_t \xrightarrow{\text{1-subG.}}$$

\Rightarrow Standard unstructured bandit environment

$$(b) A_t = \{\psi(c_t, i), i=1, 2, \dots, k\}$$

\Rightarrow contextual linear bandit.

Linear Bandit:

(see Sec 19.2 in textbook)

$$X_t = \langle \theta^*, A_t \rangle + \eta_t \quad : \text{Reward at time } t$$

$$A_t \in \mathcal{A}_t, \quad A_t \in \mathbb{R}^d$$

$$\theta^* \text{ unknown}, \quad \theta^* \in \mathbb{R}^d.$$

For this class: (constraints on action space).

a) $|\langle \theta^*, a \rangle| \leq 1 \quad , \quad \forall a \in \bigcup_{t=1}^n \mathcal{A}_t$

b) $\|a\|_2 \leq L < \infty \quad \forall a \in \bigcup_{t=1}^n \mathcal{A}_t$

Information at end of time ($t-1$):

$$H_{t-1} = \{A_1, A_1, X_1, A_2, A_2, X_2, \dots, A_{t-1}, A_{t-1}, X_{t-1}\}$$

Algorithm: Lin UCB (Linear UCB).

Note: There are several similar variants, including

- * LinRel (Linear Reinforcement Learning)
- * DFUL (Optimism in the Face of Uncertainty for Linear bandits)

Idea: Given H_{t-1} , estimate a set $C_t \subseteq \mathbb{R}^d$ such that

- ① (a) $\theta^* \in C_t$ w.h.p.
 (b) C_t gets smaller with increasing t

(ideally, $C_t = \{\theta^*\}$).

- ② Use C_t to determine an optimistic estimate of θ^* using the index for arm $a \in \mathbb{R}^d$ (note: we have an infinite number of arms in this setting) as follows:

$$UCB_t(a) = V_t(a) = \max_{\theta \in C_t} \langle \theta, a \rangle$$

- ③ Play:

$$A_t = \max_{a \in A_t} V_t(a).$$

Construction of Confidence Ellipsoid $C_t \subseteq \mathbb{R}^d$:

Recall: At begining of time t , we have observed

$$(A_1, x_1), (A_2, x_2), \dots, (A_{t-1}, x_{t-1}),$$

where $x_s = \langle \theta^*, A_s \rangle + \eta_s$

\downarrow \downarrow \downarrow \rightarrow
 $\in \mathbb{R}$ $\in \mathbb{R}^d$ $\perp\text{-subG noise}$

Approach: Use a regularized least-square estimator to estimate $\hat{\theta}_{t-1}$, and use an ellipsoid centered at $\hat{\theta}_{t-1}$, and with axes depending on the empirical covariance matrix of the estimator.

$$\hat{\theta}_{t-1} = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \left(\sum_{s=1}^{t-1} (x_s - \langle \theta, A_s \rangle)^2 + \lambda \|\theta\|_2^2 \right).$$

$\hookrightarrow \textcircled{*}$

$\lambda > 0$, regularizer. For now, assume $\lambda > 0$.

C_t construction: Solve $\hat{\theta}$:

$$\hat{\theta}_{t-1} = V_{t-1}^{-1} \sum_{s=1}^{t-1} A_s x_s$$

$\in \mathbb{R}^d$

$$V_{t-1} = \left(\gamma I + \sum_{s=1}^{t-1} A_s A_s^\top \right), \quad V_0 = \gamma I.$$

\downarrow

$d \times d$, symmetric, p.d. $\Rightarrow V_{t-1}^{-1}$ exists

$$C_t = \left\{ \theta \in \mathbb{R}^d : \|\theta - \hat{\theta}_{t-1}\|_{V_{t-1}}^2 \leq \beta_t \right\}$$

$$= \left\{ \theta \in \mathbb{R}^d : (\theta - \hat{\theta}_{t-1})^\top V_{t-1} (\theta - \hat{\theta}_{t-1}) \leq \beta_t \right\},$$

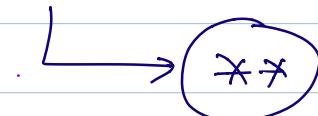
with

$\{\beta_t, t=1, 2, \dots, n\}$ an increasing sequence ($\beta_t \geq 1$):

$$\sqrt{\beta_n} = \sqrt{\gamma L} + \sqrt{2 \ln(\frac{1}{\delta}) + d \ln\left(\frac{d\gamma + nL^2}{d\gamma}\right)}$$

\downarrow time horizon \downarrow action norm-bound \downarrow prob. $\in (0, 1)$,
chosen as $\delta = \frac{1}{n}$. \downarrow dimension.

Roughly, $\beta_n \sim C' \ln(n)$



Note: Suppose $A_C = \{e_1, e_2, \dots, e_d\}$, i.e., the unit vectors in d -dimensions (essentially, a d -armed unstructured bandit, because $\langle \theta^*, e_j \rangle = \theta_j^*$).

Let $T_j(t-1)$ = number of times e_j was played until time $(t-1)$.

$$\text{Observe } A_S A_S^\top = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} (0 - 0 \Delta 0 \dots 0) = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

jth position

$$\therefore V_{t-1} = \lambda I + \begin{pmatrix} T_1(t-1) & 0 \\ T_2(t-1) & \ddots & 0 \\ \vdots & \ddots & T_d(t-1) \end{pmatrix}$$

$$\Rightarrow \hat{\theta}_{t-1} = \begin{pmatrix} \ddots & 0 \\ 0 & \frac{1}{\lambda + T_j(t-1)} \end{pmatrix} \begin{pmatrix} \ddots & \sum x_s \\ 0 & \ddots & \ddots \end{pmatrix}$$

$$= \begin{pmatrix} \ddots & \sum x_s \\ \frac{1}{\lambda + T_j(t-1)} & \ddots & \ddots \end{pmatrix}$$

(i.e., empirical mean, but with λ over-correction.)

$\{s \in J\}$ shorthand

for $\{s : A_s = e_j\}$.

Asymptotically ok as $T_j \rightarrow \infty$

$$C_t = \left\{ (\theta - \hat{\theta}_{t-1})^\top \begin{pmatrix} \ddots & & 0 \\ & \ddots & \lambda + \tau_j(t-1) \\ & & \ddots \end{pmatrix} (\theta - \hat{\theta}_{t-1}) \leq \beta_t \right\}$$

$$= \left\{ \theta \in \mathbb{R}^d : \sum_{j=1}^d (\theta_j - \hat{\theta}_{t-1,j})^2 (\lambda + \tau_j(t-1)) \leq \beta_t \right\}$$

$$UCB_t(e_j) = U_t(e_j) = \max_{\theta \in C_t} \langle \theta, e_j \rangle$$

choose

$$\theta_i = \hat{\theta}_{t-1,i} + \underbrace{\theta_j}_{\text{so terms are } d \text{ and}} \quad \forall i \neq j \quad \text{(provides more flexibility along } j^{\text{th}} \text{ coordinate in } C_t)$$

$$= \left(\hat{\theta}_{t-1,j} + \sqrt{\frac{\beta_t}{\lambda + \tau_j(t-1)}} \right).$$

$$= \frac{\sum_{s \in j} x_s}{\lambda + \tau_j(t-1)} + \sqrt{\frac{\beta_t}{\lambda + \tau_j(t-1)}} \quad \text{Scales as } \ln(t).$$

i.e., we get back the usual UCB algorithm.

In general, with more actions, arms provide richer information about all coordinates, thus providing a ellipsoid confidence region.

Regret Analysis: Take 1

Assumptions: (19.1 in textbook)

$$1. \langle \theta^*, a \rangle \leq \Delta \quad \forall a \in \bigcup_{t=1}^n A_t \quad \left. \right\}$$

$$2. \|a\|_2 \leq L \quad \forall a \in \bigcup_{t=1}^n A_t \quad \left. \right\}$$

$$3. \exists \delta \in (0, 1) \text{ s.t. with prob } \geq (1-\delta),$$

$$\forall t \in \{1, 2, \dots, n\}, \quad \theta^* \in C_t \quad \begin{matrix} \xrightarrow{\text{confidence set}} \\ \text{defined above} \end{matrix}$$

→ This is unjustified at present. C_t is centered around $\hat{\theta}_t$. To justify this assumption, we need to show $\|\theta^* - \hat{\theta}_t\|$ is "small", and that C_t is large enough to contain θ^* .

In future takes on this, this will be justified more carefully.

Thm (19.2 in textbook). Suppose assumption above holds.

Then w.p. $\geq (1-\delta)$,

$$\hat{R}_n \doteq \left(\sum_{t=1}^n \max_{\alpha_t \in A_t} \langle \theta^*, \alpha_t \rangle - \sum_{t=1}^n x_t \right)$$

$$\leq \sqrt{8n \beta_n \ln \left(\frac{\det(N_n)}{\det(V_0)} \right)}$$

$$\leq \sqrt{8n \beta_n \ln \left(\frac{d\lambda + nL^2}{d\lambda} \right)}$$

Corollary (19.3 in textbook): Suppose β_n chosen as \checkmark and $\delta = 1/n$. Then,

$$R_n \leq Cd\sqrt{n} \ln(nL),$$

$C > 0$ a fixed constant.

(19.4 in text)

Lemma: V_0 a p.d. matrix, $\{x_1, x_2, \dots, x_n\} \in \mathbb{R}^d$, with $\max_{1 \leq k \leq n} \|x_k\|_2 \leq L < \infty$. Then, with $V_t = V_0 + \sum_{s=1}^t x_s x_s^\top$,

$$\sum_{t=1}^n \left(1 + \underbrace{x_t^\top V_{t-1}^{-1} x_t}_{\|x_t\|_{V_{t-1}}^2} \right) \leq 2 \ln \left(\frac{\det(V_n)}{\det(V_0)} \right)$$

$$\leq 2d \ln \left(\frac{\text{trace}(V_0) + nL^2}{d \cdot (\det(V_0))^{1/d}} \right).$$

Pf: For any $u \geq 0$, we have $(u \wedge 1) \leq 2 \ln(u + 1)$.

Thus,

$$\sum_{t=1}^n \left(1 + x_t^\top V_{t-1}^{-1} x_t \right) \leq \sum_{t=1}^n \ln \left(1 + x_t^\top V_{t-1}^{-1} x_t \right)$$

$$\text{Now, recall } V_t = V_{t-1} + x_t x_t^\top$$

$$= V_{t-1}^{\frac{1}{2}} \left(I + V_{t-1}^{-\frac{1}{2}} x_t x_t^\top V_{t-1}^{-\frac{1}{2}} \right) V_{t-1}^{\frac{1}{2}}$$

$$\therefore \det(V_t) = \det(V_{t-1}) \det \left(I + V_{t-1}^{-\frac{1}{2}} x_t x_t^\top V_{t-1}^{-\frac{1}{2}} \right)$$

$$= \det(V_{t-1}) \cdot \det \left(I + (V_{t-1}^{-\frac{1}{2}} x_t) (V_{t-1}^{-\frac{1}{2}} x_t)^\top \right)$$

$$\text{Now } \det(I + yy^\top) = (1 + y^\top y) \quad \left(\because \text{ eigenvalues are } 1 + y^\top y, 1, 1, \dots, 1 \right)$$

$$= \det(V_{t-1}) (1 + x_t^\top V_{t-1}^{-1} x_t).$$

$$\therefore \frac{\det(V_t)}{\det(V_{t-1})} = (1 + x_t^\top V_{t-1}^{-1} x_t)$$

$$\Rightarrow \sum_{t=1}^n \ln \left(1 + x_t^\top V_{t-1}^{-1} x_t \right) = \sum_{t=1}^n \ln \left(\frac{\det(V_t)}{\det(V_{t-1})} \right)$$

$$= 2 \frac{\det(V_n)}{\det(V_0)}.$$

$$\Rightarrow \sum_{t=1}^n (1 + x_t^\top V_{t-1}^{-1} x_t) \leq 2 \frac{\det(V_n)}{\det(V_0)}.$$

Now recall V_n is p.d. with rank = d.

$$\left(\det(V_n) \right)^d = \left(\prod_{i=1}^d \lambda_i \right)^d \leq \left(\frac{1}{d} \sum_{i=1}^d \lambda_i \right)^d$$

eigenvalues
of V_n AM-GM
inequality $= \text{trace}(V_n)$

$$\leq \left(\frac{\text{trace}(V_0) + nL^2}{d} \right). \quad \square$$

PF of Thm 19.2: Recall from Assumption that $\theta^* \in C_t$
 $\forall t = 1, 2, \dots, n$. Let $A_t^* = \underset{a \in A_t}{\operatorname{argmax}} \langle \theta^*, a_t \rangle$

$$\text{and } r_t = \langle \theta^*, A_t^* - A_t \rangle$$

Recall the LinUCB algorithm: $V_t(a) = \max_{\theta \in C_t} \langle \theta, a \rangle$

$$A_t = \underset{a \in A_t}{\operatorname{argmax}} V_t(a)$$

$$\text{Let } \tilde{\theta}_t \in C_t \text{ be s.t. } \langle \tilde{\theta}_t, A_t \rangle = V_t(A_t)$$

(i.e., the upper confidence vector corresponding to the action taken).

$$\therefore \langle \theta^*, A_t^* \rangle \leq V_t(A_t^*) \leq V_t(A_t) = \langle \tilde{\theta}_t, A_t \rangle.$$

\downarrow \downarrow
 $\theta^* \in C_t$ A_t is the optimizing action
 \therefore in A_t

$$r_t = \langle \theta^*, A_t^* - A_t \rangle = \langle \theta^*, A_t^* \rangle - \langle \theta^*, A_t \rangle$$

$$\leq \langle \tilde{\theta}_t, A_t \rangle - \langle \theta^*, A_t \rangle = \langle \tilde{\theta}_t - \theta^*, A_t \rangle$$

Cauchy Schwartz inequality. $\leq \|A_t\|_{V_{t-1}^{-1}} \|\tilde{\theta}_t - \theta^*\|_{V_{t-1}}$ $\tilde{\theta}_t - \theta^*$ is symmetric, p.d. matrix

defn. of C_t $\leq 2 \|A_t\|_{V_{t-1}^{-1}} \sqrt{\beta_t}$

Further, by Assumption ①, $|\langle \varrho^*, a \rangle| \leq 1 \quad \forall a \in A_t$

$$\Rightarrow r_t \leq 2.$$

$\{\beta_t, t \geq 1\}$ increasing $\Rightarrow \beta_n \geq \max\{1, \beta_1\}.$

$$\therefore r_t \leq 2 \wedge 2\sqrt{\beta_t} \|A_t\|_{V_{t-1}^{-1}}$$

$$\leq 2\sqrt{\beta_t} (1 \wedge \|A_t\|_{V_{t-1}^{-1}})$$

$$\therefore \hat{R}_n = \sum_{t=1}^n r_t \leq \sqrt{n \sum_{t=1}^n r_t^2}$$

Cauchy-Schwartz
ineq.

$$\begin{aligned} \sum_{t=1}^n r_t &= n \cdot \frac{1}{n} \sum_{t=1}^n r_t \\ &\leq n \sqrt{\frac{1}{n} \sum_{t=1}^n r_t^2} \\ &= \sqrt{n \sum_{t=1}^n r_t^2} \end{aligned}$$

$$\leq 2\sqrt{n \beta_n \sum_{t=1}^n (1 \wedge \|A_t\|_{V_{t-1}^{-1}})}$$

Lemma

$$\leq \sqrt{8n \beta_n \ln \left(\frac{d\gamma + nL^2}{d\gamma} \right)}$$

□

Cauchy Schwartz Inequality.

$$\langle a, b \rangle = a^\top V^{-\frac{1}{2}} V^{\frac{1}{2}} b = (V^{-\frac{1}{2}} a)^\top (V^{\frac{1}{2}} b)$$

$$\leq (a^\top V^{-\frac{1}{2}} a)^{\frac{1}{2}} (b^\top V^{\frac{1}{2}} b)^{\frac{1}{2}} \\ = \|a\|_V \|b\|_V$$