Recap Stochastic Bandits

Finite arms
- Non-adaptive - uniform exploration, eps-greedy
- Adaptive - successive elimination, UCB (upper confidence bnd) sampling

Infinite arms/actions

Structured bandits - Lipschitz, Linear, GP, ...

Lipschitz bandits

$$\rightarrow |\mu(x) - \mu(x')| \leq \underbrace{L}\,\underbrace{|x-x'|} \qquad \forall x, x' \in \mathcal{Y}$$

$$\underset{\text{metric}}{D(x,x')}$$

① Fixed discretization - N bins

$$E[R(T)] = \underbrace{T\mu^*(x) - T\mu_N^*}_{\text{discretization error}} + \underbrace{T\mu_N^* - \sum_{t=1}^{T} \mu^*(x_t)}_{\text{cumulative regret from N-armed bandit}}$$
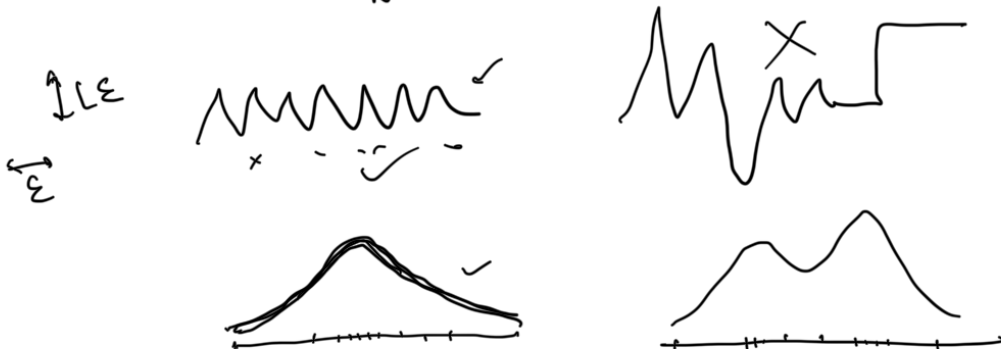
2-d

1-dim

$$\leq T\frac{L}{N} + \sqrt{NT\log T}$$

$$= O\left((L\log T)^{1/3} T^{2/3}\right)$$

Optimal in worst case for L-lipschitz rewards

Lower bound $\Omega(L^{1/3} T^{2/3})$

d-dim

$$\leq T\frac{L}{N^{1/d}} + \sqrt{NT\log T} \;\; = (L\log T)^{\frac{1}{d+2}} T^{\frac{d+1}{d+2}}$$

UB
LB.
(upto log)

② Adaptive Discretization - Zooming Algorithm

more points/arms in promising regions

active arms $S \leftarrow \phi$

For $t = 1, 2, \ldots$

    if some arm is not covered by confidence ball of active arms

    then pick any such arm & add to $S$     $x : |x - a| > \varepsilon_t(a)$
                                                        $a \in S$

    Play active arm with largest $\hat{\mu}_t(x) + 2\varepsilon_t(x)$

For any $x$    $|\hat{\mu}_t(x) - \mu(x)| \leq \underbrace{\sqrt{\frac{2 \log 1/\delta}{n_t(x)}}}_{\varepsilon_t(x)}$    $wp \geq 1 - \delta$   *     $\underline{\underline{L = 1}}$

Need to hold $\forall x$ active & all $t = 1, \ldots T$ $\leftarrow$ easy

    hard $\because$ infinitely many arms

Let $a_t$ be arm activated at time $t$. $\overbrace{\qquad\qquad}$ events are independent.

     $Pr(\text{* holds for } a_t) = \sum_x Pr(a_t = x) \cdot Pr(\text{* holds for } x)$

                              $\geq 1 - \delta$

Apply union bound $\forall a_t$ & all $t$     $\geq 1 - \delta T^2 \approx 1 - \frac{1}{T^2}$   $\delta \approx \frac{1}{T^4}$

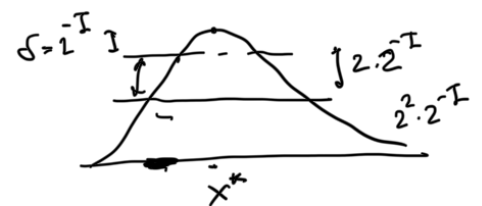Assume this high prob event from now on.

$R(T) = T\mu^* - \sum_{t=1}^{T} \mu(x_t)$

                $I = \log 1/\delta$
$\qquad\quad = T\delta + \sum_{i=1}^{I} R_i(T)$

$R_i(T) = \sum_{x : 2^{-i} \leq \Delta(x) \leq 2 \cdot 2^{-i}} n(x) \cdot \Delta(x)$

Consider active arms $\Delta(x) \leq \delta$
                       $= \mu^* - \mu(x)$

& active arms gap $\Delta(x) > \delta = 2^{-I}$

     $2^{-i} \leq \Delta(x) \leq 2 \cdot 2^{-i}$

$\delta = 2^{-I}$



We first prove the following lemma:

Lemma: $\underline{\underline{\Delta(x)}} \leq \underline{\underline{3\varepsilon_t(x)}} = O\left(\sqrt{\frac{\log T}{n_x}}\right)$   for each $x, t$.   whp

    1. $x^*$ be covered by active arm $y$

Let $x$ be chosen at time $t$.  ($x$ - active arm)

$$\mu(x) + 3\varepsilon_t(x) \geq \hat{\mu}_t(x) + 2\varepsilon_t(x) \geq \hat{\mu}_t(y) + 2\varepsilon_t(y) \geq \mu(y) + \varepsilon_t(y) \geq \mu(x^*)$$

since $y$ covers $x^*$, $\mu(y) - \mu(x^*))$
$\leq |y - x| \leq \varepsilon_t(y)$

$$\Delta(x) = \mu(x^*) - \mu(x) \leq 3\varepsilon_t(x)$$

∎

$$\Rightarrow n(x) = O\left(\frac{\log T}{\Delta^2(x)}\right)$$

$$R_i(T) = \sum_{x:\, 2^{-i} \leq \Delta(x) \leq 2 \cdot 2^{-i}} n(x) \cdot \Delta(x)$$

$$= O\left(\sum_{x:\, 2^{-i} \leq \Delta(x) \leq 2 \cdot 2^{-i}} \frac{\log T}{\Delta(x)}\right) = O\left(\sum_{x:\, 2^{-i} \leq \Delta(x) \leq 2 \cdot 2^{-i}} \frac{\log T}{2^{-i}}\right)$$

$\dfrac{2 \cdot 2^{-i}}{2^{-i}} = $ 

How many arms are activated in $\{x:\, 2^{-i} \leq \Delta(x) \leq 2 \cdot 2^{-i}\} = \Delta_i$ ?

To bound how many arms activated in $\Delta_i$, we will use the lemma above to argue that two active arms can't be too close.

Let $x, y$ are activated 4 in $\Delta_i$

$x$ activated before $y$.

When $y$ is activated, it is not covered by $x$.

$$\Rightarrow D(x,y) > \varepsilon_t(x) \geq \frac{\Delta(x)}{3} \qquad \text{from lemma,}$$

$$D(x,y) \geq \frac{1}{3}\min\, \Delta(x), \Delta(y) \geq \frac{2^{-i}}{3}$$

$$\Rightarrow R_i(T) = O\left(\frac{\log T}{2^{-i}}\, N_{\frac{2^{-i}}{3}}(x:\, 2^{-i} \leq \Delta(x) \leq 2 \cdot 2^{-i})\right)$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxx}}_{\text{Zooming dimension.}}$$

$$\underset{d \geq 0}{\inf} \; \{ N_{r/3}(\Delta_r) \leq c \cdot r^{-d} \} \; \forall r > 0$$

$$R(T) = \delta T + O\left( \frac{\log T}{\delta} \delta^{-d} \right)$$

$$\simeq O\left( T^{\frac{d+1}{d+2}} (\log T)^{\frac{1}{d+2}} \right)$$

whp

$d$ ~ zooming dim
NOT necessarily
ambient dim.

---

Linear Bandits , Gaussian bandits. NN bandits.

Concentration Bounds     Dependent data.

Independent data
Hoeffding's inequality      $X_1 \ldots X_n$ iid  mean $\mu$, $a_i \leq X_i \leq b_i$ a.s.

then    $P\left( | \frac{1}{n} \sum_{i=1}^{n} X_i - \mu | \geq \varepsilon \right) \leq e^{-2n^2 \varepsilon^2 / \sum_{i=1}^{n} (b_i - a_i)^2}$

$$e^{-2n\varepsilon^2/c^2} \;\; \Longleftarrow$$

Union bound
$$P(A \cup B) \leq P(A) + P(B)$$

Bernstein inequality      $X_1 \ldots X_n$ iid   mean $\mu$  with $E[e^{t(X_i - \mu)}] \leq e^{\frac{var(X_i) t^2 /2}{1 - b|t|}}$

then   $P\left( | \frac{1}{n} \sum_{i=1}^{n} X_i - \mu | \geq \varepsilon \right) \leq 2 e^{\frac{-n\varepsilon^2/2}{var(x) + b\varepsilon}}$      for any $t \in \left(-\frac{1}{b}, \frac{1}{b}\right), b > 0$

if $var(x)$ is small     eg. if $|X_i| \leq c$    $b = c/3$

$$\simeq e^{-n\varepsilon} \;\; \Longleftarrow$$

Martingale — seq$^n$ of random variables s.t.   $Z_1 Z_2 \ldots Z_n \ldots$

$$\forall n \quad E[Z_{n+1} | Z_1 \ldots Z_n] = Z_n \quad \Longleftarrow$$

For our purposes, martingales behave like ind. r.v..

$\rightarrow$ Azuma-Hoeffding inequality    $\{Z_i\}_{i=1}^{n}$ , $Z_1 = 0$   be a martingale

with a.s. bounded increments $|Z_i - Z_{i-1}| \le b_i$ then

$$P(Z_n \ge \epsilon) \le \exp\left\{-\frac{\epsilon^2}{2\sum_{i=1}^{n} b_i^2}\right\} \qquad \forall \epsilon > 0, n$$

Eqr. $\{X_i\}_{i=1}^{n}$ be a martingale difference seq$^n$ with $|X_i| \le b_i$ a.s.

$$P\left(\sum_{i=1}^{n} X_i \ge \epsilon\right) \le \exp\left\{-\frac{\epsilon^2}{2\sum_{i=1}^{n} b_i^2}\right\}$$

$$Z_n = \sum_{i=1}^{n} X_i \qquad E[X_{n+1}] = 0 . \qquad E\left[\sum_{i=1}^{n+1} X_i \mid X_1, \dots X_n\right] = \sum_{i=1}^{n} X_i \Leftarrow$$

$$\overset{1)}{=} E[X_{n+1}] + \sum_{i=1}^{n} X_i$$

$$\underbrace{\qquad}_{=0}$$

$\Rightarrow$ Bernstein inequality for Martingales (Freedman's inequality)

$\{X_i\}_{i=1}^{n}$ be a martingale difference seq$^i$ with $|X_i| \le b_i$ a.s.

$$P\left(\sum_{i=1}^{n} X_i \ge t\right) \le \exp\left\{\frac{-\epsilon^2}{\underbrace{\sum_{t=1}^{n} E[X_i^2 \mid X_{1:i-1}]}_{\text{var}} + \frac{b\epsilon}{3}}\right\} \approx e^{-\epsilon}$$
$$\text{if var}$$
$$\text{is small}$$

---

$$\overset{*}{\mu(x)} = \theta^{*T} x \qquad \text{Linear reward.} \qquad \mu(x) = \theta^T x$$

whp $\rightarrow \hat{\mu}_t(x) - \mu(x) \approx \theta^* - \hat{\theta}_t \qquad \theta, \theta^* - d\text{-dim}$

$$\Rightarrow \|\theta^* - \hat{\theta}_t\|_{V_t}^2 \le \beta_t \approx \sigma d$$

$$\rightarrow \hat{\theta}_t = \left(\sum_{s=1}^{t} x_s x_s^T\right)^{-1} x_s^T y_s + \lambda I \qquad V_t = \sum_{s=1}^{t} x_s x_s^T$$