

Lecture 08: Two and Multi-Armed Bandit

Lecturer: Yuan Zhou

Scribe: Zexing Xu, Erik Pan

1 Recap

$\exists A : (\epsilon, \delta) - PAC$ for the better arm, sample complexity $\leq T \Rightarrow \Delta(D_1, D_2) \geq 1 - 2\delta$, where

$$D_1 = B_p^{\otimes T} \otimes B_{p+\epsilon}^{\otimes T}$$

$$D_2 = B_{p+\epsilon}^{\otimes T} \otimes B_p^{\otimes T}$$

, where Δ is total variation distance, D_1, D_2 are two different distributions defined on the same probability space and B_p is Bernoulli distribution with parameter p .

Pinsker's inequality (proved in Lecture 07) stated

$$\textcircled{1} : \Delta(P, Q) \leq \sqrt{D_{KL}(P||Q)}$$

$$\textcircled{2} : \Delta(P, Q) \leq 1 - \frac{1}{2} \exp(-D_{KL}(P||Q))$$

If D_{KL} is small (i.e. $\approx \epsilon$) then

$$\textcircled{1} : \Delta(P, Q) \leq \sqrt{\epsilon}$$

$$\textcircled{2} : \Delta(P, Q) \leq 1 - \frac{1}{2} e^{-\epsilon}$$

So $\textcircled{1}$ gives a better bound.

While if D_{KL} is large (i.e. ≈ 100) then

$$\textcircled{1} : \textit{trivial}$$

$$\textcircled{2} : \Delta(P, Q) \leq 1 - \frac{1}{2} e^{-100} < 1$$

So $\textcircled{2}$ gives a better bound.

2 Two-Armed Bandit

As discussed in Lecture 7, K-L Divergence has some limitations, but why it is still frequently used in measuring the variational difference? The answer lies in the following promising properties of K-L Divergence.

Fact (Additivity): Suppose $P(x, y) = P_1(x)P_2(y), Q(x, y) = Q_1(x)Q_2(y)$, then :

$$\begin{aligned}
D_{KL}(P, Q) &= - \sum_{x,y} P(x,y) \cdot \ln\left(\frac{Q(x,y)}{P(x,y)}\right) \\
&= - \sum_{x,y} P_1(x)P_2(y) [\ln\left(\frac{Q_1(x)}{P_1(x)}\right) + \ln\left(\frac{Q_2(y)}{P_2(y)}\right)] \\
&= - \sum_x P_1(x) [\ln\left(\frac{Q_1(x)}{P_1(x)}\right)] - \sum_y P_2(y) [\ln\left(\frac{Q_2(y)}{P_2(y)}\right)] \\
&= D_{KL}(P_1||Q_1) + D_{KL}(P_2||Q_2)
\end{aligned}$$

Since D_{KL} is not symmetric, we can't combine the 2 terms, so we have :

$$D_{KL}(D_1||D_2) = T \cdot D_{KL}(B_p||B_{p+\epsilon}) + T \cdot D_{KL}(B_{p+\epsilon}||B_p)$$

Since :

$$\begin{aligned}
D_{KL}(B_p||B_{p+\epsilon}) &= -p \cdot \ln\left(1 + \frac{\epsilon}{p}\right) - (1-p) \cdot \ln\left(1 + \frac{\epsilon}{1-p}\right) \\
&= -p\left(\frac{\epsilon}{p} - \frac{1}{2}\frac{\epsilon^2}{p^2} + \mathcal{O}\left(\frac{\epsilon^3}{p^3}\right)\right) - (1-p)\left(-\frac{\epsilon}{1-p} - \frac{1}{2}\frac{\epsilon^2}{(1-p)^2} + \mathcal{O}\left(\frac{\epsilon^3}{(1-p)^3}\right)\right) \quad (\text{Taylor Expansion}) \\
&= \epsilon^2 \cdot \left(\frac{1}{2p} + \frac{1}{2 \cdot (1-p)}\right) + \mathcal{O}(\epsilon^3) \cdot \left(\frac{1}{p^3} + \frac{1}{(1-p)^3}\right) \\
&= \mathcal{O}(\epsilon^3) \quad (\text{as long as } p \text{ is away from } 0, 1 \text{ (e.g. } p \in [0.1, 0.9]) \text{ and } \epsilon \text{ is small)}
\end{aligned}$$

So :

$$D_{KL}(D_1||D_2) = \mathcal{O}(\epsilon^2 \cdot T)$$

So by Pinsker's inequality ① :

$$\Delta(D_1, D_2) \leq \mathcal{O}(\epsilon \cdot \sqrt{T})$$

when $\delta > 0.01$,

$$\left. \begin{aligned} \Delta(D_1, D_2) &\geq 0.9 \\ \Delta(D_1, D_2) &\leq \mathcal{O}(\epsilon \cdot \sqrt{T}) \end{aligned} \right\} \Rightarrow \epsilon \cdot \sqrt{T} \geq \Omega(1) \Rightarrow T \gtrsim \frac{1}{\epsilon^2}$$

For constant δ , the algorithm give us $T \lesssim \frac{1}{\epsilon^2} \log(\frac{1}{\delta})$, and δ is constant, the two bounds matches.

when $\delta \ll 0.01$ (i.e. δ is very small), by Pinsker's inequality ② :

$$\left. \begin{aligned} \Delta(D_1, D_2) &\leq 1 - \frac{1}{2}e^{-\Theta(T\epsilon^2)} \\ \Delta(D_1, D_2) &\geq 1 - 2\delta \end{aligned} \right\} \Rightarrow e^{-\Theta(T\epsilon^2)} \leq 4\delta \Rightarrow -\Theta(T\epsilon^2) \leq \ln(4\delta) \Rightarrow T\epsilon^2 \gtrsim \ln\left(\frac{1}{\delta}\right) \Rightarrow T \gtrsim \epsilon^{-2} \ln\left(\frac{1}{\delta}\right)$$

3 Multi-Armed Bandit

For any bandit instance, let $I = (D_1, D_2, \dots, D_n)$ be the reward distribution for each arm i , and (possibly randomize) policy π , running π on I result the history $H = (a_1, r_1, a_2, r_2, \dots, a_{|H|}, r_{|H|})$ where $|H|$ is the length of H . Let $P_{\pi, I}$ be the distribution of H .

Lemma 1. Let $I = (D_1, D_2, \dots, D_n), I' = (D'_1, D'_2, \dots, D'_n)$ for any π :

$$D_{KL}(P_{\pi, I} || P_{\pi, I'}) = \sum_{i=1}^n \mathbb{E}_{H \sim P_{\pi, I}} [T_i(H)] \cdot D_{KL}(D_i || D'_i) \quad (\text{where } T_i[H] \text{ is how many times arm } i \text{ is played in } H)$$

Lemma 1 describes the relationship between the whole distribution and the original individual distributions. Intuitively, if we run K-L divergence on similar distance, K-L does not diverge much.

Proof.

$$\begin{aligned} P_{\pi, I}(H) &= P_{\pi}(a_1) \cdot P_{D_{a_1}}(r_1) \cdot P_{\pi}(a_2) \cdot P_{D_{a_2}}(r_2) \cdot P_{\pi}(a_3) \cdot P_{D_{a_3}}(r_3) \cdot \dots \\ &= \prod_{t=1}^{|H|} P_{\pi}(a_t | a_1, r_1, a_2, r_2, \dots, a_{t-1}, r_{t-1}) \cdot P_{D_{a_t}}(r_t) \end{aligned}$$

Similarly, we have

$$P_{\pi, I'}(H) = \prod_{t=1}^{|H|} P_{\pi}(a_t | a_1, r_1, a_2, r_2, \dots, a_{t-1}, r_{t-1}) \cdot P_{D'_{a_t}}(r_t)$$

For $P_{\pi, I'}(H) \neq 0$, we have

$$\ln\left(\frac{P_{\pi, I}(H)}{P_{\pi, I'}(H)}\right) = \sum_{t=1}^{|H|} \ln\left(\frac{P_{D_{a_t}}(r_t)}{P_{D'_{a_t}}(r_t)}\right)$$

Then it follows that

$$D_{KL}(P_{\pi, I}(H) || P_{\pi, I'}(H)) = \mathbb{E}_{H \sim P_{\pi, I}} \ln\left(\frac{P_{\pi, I}(H)}{P_{\pi, I'}(H)}\right) = \mathbb{E}_{H \sim P_{\pi, I}} \sum_{t=1}^{|H|} \ln\left(\frac{P_{D_{a_t}}(r_t)}{P_{D'_{a_t}}(r_t)}\right)$$

$$\begin{aligned} \sum_{t=1}^{\infty} \mathbb{E}_{H \sim P_{\pi, I}} \mathbb{I}[t \leq |H|] \ln\left(\frac{P_{D_{a_t}}(r_t)}{P_{D'_{a_t}}(r_t)}\right) &= \sum_{t=1}^{\infty} \mathbb{E}_{H \sim P_{\pi, I}} \mathbb{I}[t \leq |H|] \sum_{i=1}^n \mathbb{I}[a_t = i] \ln\left(\frac{D_i(r_t)}{D'_i(r_t)}\right) \quad (\text{where } \mathbb{I}[\cdot] \text{ is indicator function}) \\ &= \sum_{t=1}^{\infty} \sum_{i=1}^n \mathbb{E}_{(a_1, t_1, \dots, a_{t-1}, r_{t-1}) \sim P_{\pi, I}} \mathbb{I}[a_t \text{ exist} \wedge a_t = i] \mathbb{E}_{r_t \sim D_i} \ln\left(\frac{D_i(r_t)}{D'_i(r_t)}\right) \\ &= \sum_{i=1}^n \mathbb{E}_{H \sim P_{\pi, I}} [T_i(H)] - D_{KL}(D_i || D'_i) \end{aligned}$$

Note: Since a_t and r_t is independent of H , we do not need conditional expectation for the term $\mathbb{E}_{r_t \sim D_i} \ln\left(\frac{D_i(r_t)}{D'_i(r_t)}\right)$. \square

Theorem 2 (minimax regret lower bound). $\forall n \geq 2, T \geq n$, policy π , $\exists I$ such that $R_T^{\pi, I} \gtrsim \sqrt{nT}$

Proof. Let $I = (B_{\frac{1}{2}+\Delta}, B_{\frac{1}{2}}, B_{\frac{1}{2}}, \dots, B_{\frac{1}{2}})$, let $Z = \operatorname{argmin}_{i \in \{2, 3, \dots, n\}} \mathbb{E}_{H \sim P_{\pi, I}} [T_i(H)]$.

Set $I' = (B_{\frac{1}{2}+\Delta}, B_{\frac{1}{2}}, B_{\frac{1}{2}}, \dots, B_{\frac{1}{2}}, B_{\frac{1}{2}+2\Delta}, B_{\frac{1}{2}}, B_{\frac{1}{2}}, \dots, B_{\frac{1}{2}})$. Typically, we set the reward distribution for z th arm as $B_{\frac{1}{2}+2\Delta}$ to have slightly better performance than our original best arm1, but this arm may not receive

much attention.

Let A be the event $\{T_1(H) \leq \frac{T}{2}\}$ (Play 1st arm less than 1/2 of the time), clearly:

$$\begin{cases} R_T^{\pi,I} |_A \geq \frac{T}{2} \Delta \\ R_T^{\pi,I'} |_{\bar{A}} \geq \frac{T}{2} \Delta \end{cases}$$

$$\begin{aligned} R_T^{\pi,I} + R_T^{\pi,I'} &\geq \frac{T\Delta}{2} \cdot \Pr_{\pi,I}[A] + \Pr_{\pi,I'}[\bar{A}] = (\Pr_{\pi,I}[A] + (1 - \Pr_{\pi,I'}[A])) = \frac{T\Delta}{2} \cdot (1 - |\Pr_{\pi,I}[A] - \Pr_{\pi,I'}[A]|) \\ &\geq \frac{T\Delta}{2} (1 - \delta(P_{\pi,I}, P_{\pi,I'})) \\ &\geq \frac{T\Delta}{2} (1 - \sqrt{D_{KL}(P_{\pi,I} || P_{\pi,I'})}) \text{ (By Pinsker's Inequality ①)} \end{aligned}$$

Since

$$D_{KL}(P_{\pi,I} || P_{\pi,I'}) = \mathbb{E}_{\pi,I}[T_z] D_{KL}(D_Z, D'_Z) \leq \frac{T}{n-1} (2\Delta^2) \lesssim \frac{T}{n-1} \Delta^2$$

So

$$D_{KL}(P_{\pi,I} || P_{\pi,I'}) \geq \frac{T\delta}{2} (1 - \mathcal{O}(\Delta \sqrt{\frac{T}{n}}))$$

Pick $\delta = c \cdot \sqrt{\frac{n}{T}}$, where c is any constant such that $(1 - \mathcal{O}(\Delta \sqrt{\frac{T}{n}})) < \frac{1}{2}$

$$D_{KL}(P_{\pi,I} || P_{\pi,I'}) \geq \frac{T \cdot c \sqrt{\frac{n}{T}}}{4} \gtrsim \sqrt{nT}$$

□

Similarly, we can proof PAC has a lower bound of $n\epsilon^{-2} \ln(\frac{1}{\delta})$