

Lecture 8: Adversarial Bandits: EXP3 Algorithm

Lecturer: M. K. Hanawal

Scribes: Saumya Suri

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

8.1 Adversarial Bandit Setting

The hypothesis we used earlier will now be called Actions/Arms. Let there be K actions and let in each round the player selects one out of the K actions and incurs a loss. Unlike the full information setting, where the player had information about the losses due to every possible hypothesis, in the Adversarial Bandit setting, the player has the information of only the action he played in a round. Thus in each round he gets only $1/K$ of the full information. Hence we can say that T rounds in a full information setting is equivalent to KT rounds of the Adversarial bandit setting.

Let $I_t \in [K]$ be the arm played by the player in round t and $l_{i,t}$ be the loss of arm i in round t . Thus, at each time step $t = 1, 2, \dots$, the adversary assigns to each arm $i = 1, \dots, K$ a loss $l_{i,t}$. We assume that the losses are bounded such that $l_{i,t} \in [0, 1] \forall t$.

Algorithm 1 Template of the adversarial bandit setting

- 1: **for** $t = 1, 2, \dots, T$ **do**
 - 2: Adversary selects $l_{i,t}$, $i \in [K]$
 - 3: Player selects an arm $I_t \in [K]$
 - 4: Player observes only $l_{I_t,t}$ i.e, loss corresponding to only that arm
-

For example,

If suppose the adversary has assigned losses $l_{1,t}, l_{2,t}, \dots, l_{10,t}$ for 10 actions in round t and let say the player plays action 5 in round t , then he would be able to observe only $l_{5,t}$ in that round.

Regret of the algorithm that picks I_t in round t is:

$$\tilde{R}_T = \sum_{t=1}^T l_{I_t,t} - \min_{k=1,2,\dots,K} \sum_{t=1}^T l_{k,t} \quad (8.1)$$

where, $\sum_{t=1}^T l_{I_t,t}$ is the total loss incurred by the player for playing I_t in T rounds and $\min_{k=1,2,\dots,K} \sum_{t=1}^T l_{k,t}$ is the minimum loss he would have incurred for playing one particular arm for T rounds.

The randomness of the losses and the arm chosen by the player makes \tilde{R}_T a random quantity. Taking Expectation on the sides in (8.1), we get,

$$\begin{aligned} R_T = E(\tilde{R}_T) &= E\left(\sum_{t=1}^T l_{I_t,t}\right) - E\left(\min_{k=1,2,\dots,K} \sum_{t=1}^T l_{k,t}\right) \\ &\geq E\left(\sum_{t=1}^T l_{I_t,t}\right) - \min_{k=1,2,\dots,K} E\left(\sum_{t=1}^T l_{k,t}\right) \end{aligned}$$

The term $E\left(\sum_{t=1}^T l_{I_t,t}\right) - \min_{k=1,2,\dots,K} E\left(\sum_{t=1}^T l_{k,t}\right)$ is known as **Pseudo Regret** and is denoted by \bar{R}_T . Thus we have,

$$R_T \geq \bar{R}_T$$

We use Pseudo regret instead of the usual Expected regret(R_T) since $E\left(\min_{k=1,2,\dots,K} \sum_{t=1}^T l_{k,t}\right)$ is difficult to compute as compared to $\min_{k=1,2,\dots,K} E\left(\sum_{t=1}^T l_{k,t}\right)$.

Oblivious Adversary: If before setting the loss value, the adversary does not know which arm that will be chosen by the player then the adversary is known as Oblivious adversary. For oblivious adversary, $R_T = \bar{R}_T$.

8.1.1 Pseudo Regret bounds

EXP3 (Exponential weights for Exploration and Exploitation) Algorithm

The parameters for the Exp3 algorithm are,

- A non increasing sequence of real numbers $(n_t)_{t \in N}$
- P_1 , a uniform distribution over $[K]$

Algorithm 2 EXP3 Algorithm

- 1: **Input:** A non increasing sequence of real numbers $(n_t)_{t \in N}$
 - 2: **for** $t = 1, 2, \dots, T$ **do**
 - 3: Draw an arm I_t from the probability distribution P_t
 - 4: Player selects an arm $I_t \in [K]$
 - 5: **for** Each arm $i = 1, 2, \dots, K$ **do**
 - 6: compute the estimated loss, $\tilde{l}_{i,t} = \frac{l_{i,t}}{P_{i,t}} \mathbb{1}_{\{I_t=i\}}$
 - 7: update the cumulative loss, $\tilde{L}_{i,t} = \tilde{L}_{i,t-1} + \tilde{l}_{i,t}$
 - 8: Compute the new probability distribution:

$$P_{t+1} = (P_{1,t+1}, P_{2,t+1}, \dots, P_{K,t+1}) \text{ where } P_{i,t+1} = \frac{\exp(\eta_t \tilde{L}_{i,t})}{\sum_{k=1,\dots,K} \exp(-\eta_t \tilde{L}_{k,t})}$$
-

The player starts with a uniform distribution P_1 and then in each round pick an arm according to the distribution P_t . We let that arm be I_t and for all arms, calculate the estimated loss,

$$\tilde{l}_{i,t} = \frac{l_{i,t}}{P_{i,t}} \mathbb{1}_{\{I_t=i\}} \quad (8.2)$$

where, $P_{i,t}$ is the probability of choosing arm i in round t .

We can further show that $\tilde{l}_{i,t}$ is an unbiased Estimator (w.r.t I_t) of $l_{i,t}$ i.e,

$$E(\tilde{l}_{i,t}) = l_{i,t}$$

since,

$$\sum_{k=1, \dots, K} P_{k,t} \frac{l_{k,t}}{P_{k,t}} \mathbb{1}_{\{k=i\}} = \sum_{k=1, \dots, K} l_{k,t} \mathbb{1}_{\{k=i\}} = l_{i,t}$$

Thus, though only the loss of the played arm is observed, we build an unbiased estimator for the loss of any other arm as well.

Theorem 8.1 Pseudo Regret of Exp3:

If Exp3 is run with $\eta_t = \eta = \sqrt{\frac{2 \ln K}{TK}}$ then the Pseudo regret, i.e, \bar{R}_T is bounded above by $\sqrt{2KT \ln K}$, i.e,

$$\bar{R}_T \leq \sqrt{2KT \ln K}$$

Moreover, if Exp3 is run with $\eta_t = \eta = \sqrt{\frac{\ln K}{tK}}$ then the Pseudo regret, i.e, \bar{R}_T is bounded above by $2\sqrt{KT \ln K}$ i.e,

$$\bar{R}_T \leq 2\sqrt{KT \ln K}$$