

Lecture 13 — September 16

Lecturer: Aditya Gopalan

Scribe: Sayak Ray Chowdhury

13.1 Recall

In the previous lecture, we have seen EXP3 algorithm regarding bandit model of learning i.e on-line learning with partial information. The original work by Auer et. al. [1] considers a slightly different version of what we have seen. Here we present the algorithm given in [1] and denote it by EXP3-ORIG.

EXP3-ORIG

Parameter: $\eta \in [0, 1]$

Initialize: $p_1 = \text{Uniform}\{1, 2, \dots, N\}$; $\tilde{L}_0 = \mathbf{0} \in \mathbb{R}^N$

At each time $t = 1, 2, 3, \dots, T$

1. Sample $I_t \sim p_t$, where $p_t \equiv (p_{i,t})_{i=1}^N$

For each $i = 1, 2, \dots, N$

2. $\tilde{l}_{i,t} := \frac{l_{i,t}}{p_{i,t}} \mathbb{1}\{I_t = i\}$, where $l_t \equiv (l_{i,t})_{i=1}^N$, $\tilde{l}_t \equiv (\tilde{l}_{i,t})_{i=1}^N$

3. $\tilde{L}_{i,t} := \tilde{L}_{i,t-1} + \tilde{l}_{i,t}$, where $\tilde{L}_t \equiv (\tilde{L}_{i,t})_{i=1}^N$

4. $p_{i,t+1} := \frac{(1 - \eta) \exp(-\eta \tilde{L}_{i,t})}{\sum_{j=1}^N \exp(-\eta \tilde{L}_{j,t})} + \frac{\eta}{N}$

From the previous lecture, we know $\mathbb{E}[\text{Regret}_T^{\text{EXP3}}] \leq O(\sqrt{NT \log N})$. This **optimal** bound holds also for EXP3-ORIG. In this lecture, we will :

- Show an EXP3 like algorithm that enjoys regret bound with high probability (WHP).
- Give a lower bound on regret (MINIMAX regret) across all bandit algorithms.

13.2 Modification of EXP3 to get WHP regret

To motivate this modification, first we argue about a technical issue (though very loose) with EXP3-ORIG and equivalently with EXP3 : Variance of estimated losses $\tilde{L}_{i,t}$, which are unbiased estimates of $L_{i,t}$ can be very large for both the algorithms. First, lets see why it is so.

Recall the definition of $\tilde{l}_{i,t}$. We have,

$$\mathbb{E}[\tilde{l}_{i,t}^2 \mid \mathcal{F}_{t-1}] = p_{i,t} \frac{l_{i,t}^2}{p_{i,t}^2} = \frac{l_{i,t}^2}{p_{i,t}} \approx O(1/p_{i,t})$$

Now, in EXP3-ORIG, optimal $\eta \approx \frac{1}{\sqrt{T}}$ and $p_{i,t} \geq \frac{\eta}{N} \approx \frac{1}{N\sqrt{T}}$

So we have, $\text{Var}[\tilde{l}_{i,t} | \mathcal{F}_{t-1}] \approx O(\sqrt{T})$

Hence, $\text{Var}[\tilde{L}_{i,t} | \mathcal{F}_{t-1}] \approx O(T^{3/2})$, which gets very large as T becomes larger.

Similarly, for EXP3 also $p_{i,t}$ can be very small, thus making the variance very large. To overcome this issue, we make two key tweaks in EXP3-ORIG :

(i) Let us consider **rewards** or gains instead of losses, i.e.

gains $g_{i,t} := 1 - l_{i,t}$; $g_{i,t} \in [0, 1]$ and

gain estimates $\tilde{g}_{i,t} := \frac{g_{i,t}}{p_{i,t}} \mathbb{1}\{I_t = i\}$. Note, $\tilde{g}_{i,t} \neq 1 - \tilde{l}_{i,t}$

(ii) Control variance of (gain) estimates by adding a stabilization term (β) :

$$\begin{aligned} g'_{i,t} &:= \frac{g_{i,t} + \beta}{p_{i,t}}, \text{ if } I_t = i \\ &:= \frac{\beta}{p_{i,t}}, \text{ if } I_t \neq i \\ &= \tilde{g}_{i,t} + \frac{\beta}{p_{i,t}} \end{aligned}$$

The underlying idea behind these tweaks is to ensure that $G'_{i,t}$ is an **upper confidence bound** for

$G_{i,t}$, where $G_{i,t} = \sum_{s=1}^t g_{i,s}$ and $G'_{i,t} = \sum_{s=1}^t g'_{i,s}$

Now, we present the WHP version of the EXP3-ORIG algorithm and call it as EXP3.P as given in [1].

EXP3.P

Parameters: $\beta, \gamma, \eta \in [0, 1]$

Initialize: $p_1 = \text{Uniform}\{1, 2, \dots, N\}$; $G'_{i,0} = 0, \forall i \in [N]$

At each time $t = 1, 2, 3, \dots, T$

1. Sample $I_t \sim p_t$

For each $i = 1, 2, \dots, N$

2. $g'_{i,t} := \frac{g_{i,t} \mathbb{1}\{I_t = i\} + \beta}{p_{i,t}}$

3. $G'_{i,t} := G'_{i,t-1} + g'_{i,t}$

4. $p_{i,t+1} := \frac{(1 - \gamma) \exp(\eta G'_{i,t})}{\sum_{j=1}^N \exp(\eta G'_{j,t})} + \frac{\gamma}{N}$

Theorem 13.1. [Regret bound for EXP3.P]

For any $\delta \in (0, 1)$ with Probability $\geq (1 - \delta)$, $\text{Regret}_T^{\text{EXP3.P}} \leq 5.15 \sqrt{NT \log(N/\delta)}$

Now choosing δ to be small enough, the bound can be satisfied with high probability.

Proof: Before proving the theorem, first consider the following lemma :

Lemma 13.2. [Upper Confidence property of g']

For $\beta < 1$, with probability $\geq (1 - \delta)$, $\sum_{t=1}^T g_{i,t} \leq \sum_{t=1}^T g'_{i,t} + \frac{\log(1/\delta)}{\beta}$, $\forall i \in [N]$

Proof: Let, $\mathcal{F}_{t-1} = \sigma - alg(I_1, I_2, \dots, I_{t-1}, g(I_1, 1), g(I_2, 2), \dots, g(I_{t-1}, t-1))$, where $g(i, t)$ denote $g_{i,t}$ and let $\mathbb{E}_t[\cdot]$ denote $\mathbb{E}[\cdot | \mathcal{F}_{t-1}]$. Now,

$$\begin{aligned} & \mathbb{E}_t[\exp(\beta g_{i,t} - \beta g'_{i,t})] \\ &= \mathbb{E}_t[\exp(\beta g_{i,t} - \frac{\beta g_{i,t}}{p_{i,t}} \mathbb{1}\{I_t = i\}) \exp(\frac{-\beta^2}{p_{i,t}})] \quad [\text{from definition of } g'_{i,t}] \end{aligned}$$

$$= \mathbb{E}_t[\exp(\beta(g_{i,t} - \tilde{g}_{i,t})) \exp(\frac{-\beta^2}{p_{i,t}})] \quad [\text{from definition of } \tilde{g}_{i,t}]$$

$$= \mathbb{E}_t[1 + \beta(g_{i,t} - \tilde{g}_{i,t}) + (\beta(g_{i,t} - \tilde{g}_{i,t}))^2] \exp(\frac{-\beta^2}{p_{i,t}})$$

[using, $e^x = 1 + x + x^2$, $\forall x \leq 1$, where $x = \beta(g_{i,t} - \tilde{g}_{i,t}) \leq \beta g_{i,t} \leq 1$ and using the fact that $p_{i,t}$ is measurable w.r.t \mathcal{F}_{t-1}]

$$= (1 + \beta^2 \text{Var}_t[\tilde{g}_{i,t}]) \exp(\frac{-\beta^2}{p_{i,t}}) \quad [\text{as, } \tilde{g} \text{ is an unbiased estimator of } g]$$

$$= (1 + \beta^2 \frac{g_{i,t}^2}{p_{i,t}}) \exp(\frac{-\beta^2}{p_{i,t}})$$

$$\leq \exp(\frac{\beta^2}{p_{i,t}}(g_{i,t}^2 - 1)) \quad [\text{using, } 1 + x \leq e^x, \forall x \geq 0]$$

$$\leq 1$$

Multiplying over $t = 1, 2, \dots, T$ we get, $\mathbb{E}[\exp(\beta G_{i,t} - \beta G'_{i,t})] \leq 1$

Now **Markov's inequality** gives,

$$\text{Prob}[\beta G_{i,t} - \beta G'_{i,t} \geq \log(1/\delta)] \leq \frac{1}{1/\delta} = \delta$$

Thus, with probability $\geq (1 - \delta)$, $G_{i,t} \leq G'_{i,t} + \frac{\log(1/\delta)}{\beta}$

This is true for each $i = 1, 2, \dots, N$. □

Now, it remains to prove theorem 13.1.

Proof of theorem: We will consider the **potential function** approach like exponential weights forecaster as discussed earlier. Let,

$$W_t := \sum_{i=1}^N w_{i,t} := \sum_{i=1}^N \exp(\eta G'_{i,t}).$$

$$\text{So, } \log\left(\frac{W_T}{W_0}\right) = \log\left(\frac{\sum_{i=1}^N \exp(\eta G'_{i,t})}{N}\right) \geq \eta \max_{i=1}^N G'_{i,t} - \log N \quad \text{----- (*)}$$

On the otherhand,

$$\log\left(\frac{W_t}{W_{t-1}}\right) = \log\left(\sum_{i=1}^N \left(\frac{w_{i,t-1}}{W_{t-1}}\right) \exp(\eta g'_{i,t})\right)$$

$$\begin{aligned}
&= \log\left(\sum_{i=1}^N \left(\frac{p_{i,t} - \gamma/N}{1 - \gamma}\right) \exp(\eta g'_{i,t})\right) \quad [\text{from definiton of EXP3.P update}] \\
&= \log\left(\sum_{i=1}^N \left(\frac{p_{i,t} - \gamma/N}{1 - \gamma}\right) (1 + \eta g'_{i,t} + \eta^2 g'_{i,t}{}^2)\right) \quad [\text{setting } \eta \text{ s.t. } \eta g' \leq 1 \text{ and using} \\
&\hspace{15em} e^x \leq 1 + x + x^2, \forall x \leq 1] \\
&\leq \log\left(1 + \frac{\eta}{1 - \gamma} \sum_i p_{i,t} g'_{i,t} + \frac{\eta^2}{1 - \gamma} \sum_i p_{i,t} g'_{i,t}{}^2\right)
\end{aligned}$$

$$\begin{aligned}
\text{Now, } \sum_i p_{i,t} g'_{i,t}{}^2 &= \sum_i p_{i,t} g'_{i,t} \left(\frac{g_{i,t} \mathbb{1}\{I_t = i\} + \beta}{p_{i,t}}\right) \\
&= g'(I_t, t) g_{I_t, t} + \beta \sum_i g'_{i,t} \\
&\leq (1 + \beta) \sum_i g'_{i,t}
\end{aligned}$$

$$\text{And } \sum_i p_{i,t} g'_{i,t} = \sum_i p_{i,t} \left(\frac{g_{i,t} \mathbb{1}\{I_t = i\} + \beta}{p_{i,t}}\right) = g_{I_t, t} + N\beta$$

$$\text{Hence, } \log\left(\frac{W_t}{W_{t-1}}\right) \leq \frac{\eta}{1 - \gamma} (g_{I_t, t} + N\beta) + \frac{\eta^2}{1 - \gamma} (1 + \beta) \sum_i g'_{i,t} \quad [\text{as } \log(1 + x) \leq x, \forall x]$$

Summing over $t = 1, 2, \dots, T$ we get,

$$\log\left(\frac{W_T}{W_0}\right) \leq \frac{\eta}{1 - \gamma} G_T^{EXP3.P} + \frac{\eta N \beta T}{1 - \gamma} + \frac{\eta^2}{1 - \gamma} (1 + \beta) \sum_{i=1}^N G'_{i,T} \quad \text{----- (**)}$$

$$\text{where, } G_T^{EXP3.P} = \sum_{t=1}^T g_{I_t, t}$$

Putting (*) and (**) together we get,

$$\eta \max_i G'_{i,T} - \log N \leq \frac{\eta}{1 - \gamma} G_T^{EXP3.P} + \frac{\eta N \beta T}{1 - \gamma} + \frac{\eta^2}{1 - \gamma} (1 + \beta) \sum_{i=1}^N G'_{i,T}$$

$$\Rightarrow G_T^{EXP3.P} - (1 - \gamma) G'_{max} \geq \frac{-(1 - \gamma)}{\eta} \log N - N\beta T - \eta(1 + \beta) N G'_{max}, \text{ where } G'_{max} = \max_i G'_{i,T}$$

$$\Rightarrow G_T^{EXP3.P} \geq \frac{-\log N}{\eta} - N\beta T + (1 - \gamma - \eta(1 + \beta)N) G'_{max} \quad [\text{as, } 1 - \gamma \leq 1]$$

Now, from lemma 13.2 we know,

$$\text{w.p. } \geq (1 - \delta), G_{i,T} \leq G'_{i,T} + \frac{\log(1/\delta)}{\beta}, \forall i \in [N].$$

Applying union bound,

$$\text{w.p. } \geq (1 - \delta), \max_i G_{i,T} \leq \max_i G'_{i,T} + \frac{\log(N/\delta)}{\beta}$$

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

$$G_T^{EXP3.P} \geq \frac{-\log N}{\eta} - N\beta T + (1 - \gamma - \eta(1 + \beta)N) \max_i G_{i,T} - \frac{\log(N/\delta)}{\beta} (1 - \gamma - \eta(1 + \beta)N)$$

[Choosing $\gamma, \beta, \eta \in [0, 1]$ s.t. $0 \leq 1 - \gamma - \eta(1 + \beta)N \leq 1$]

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

$$G_{i^*, T} - G_T^{EXP3.P} \leq \frac{\log N}{\eta} + N\beta T + \frac{\log(N/\delta)}{\beta} + \gamma G_{i^*, T} + \eta(1 + \beta) G_{i^*, T}$$

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

$Regret_T^{EXP3.P} \leq \frac{\log N}{\eta} + N\beta T + \frac{\log(N/\delta)}{\beta} + \gamma T + (1 + \beta)\eta T$, where $G_{i^*,T} = \max_i G_{i,T} \leq T$

Now, we can optimize β, η, γ to get: $Regret_T^{EXP3.P} \leq 5.15\sqrt{TN\log(N/\delta)}$

Here optimal parameters are : $\beta = \sqrt{\frac{\log(N/\delta)}{T}}$, $\eta = 0.95\sqrt{\frac{\log N}{NT}}$, $\gamma = 1.05\sqrt{\frac{N\log N}{T}}$

□

13.3 Minimax lower bound across all bandit algorithms

Here, we will work in **rewards setting** and use the same notations as in the previous section.

Theorem 13.3. [Regret lower bound : bandits]

$$\inf \sup_{i=1}^N E\left[\sum_{t=1}^T g_{i,t}\right] - E\left[\sum_{t=1}^T g_{I_t,t}\right] \geq 1/20\sqrt{TN}$$

where, infimum is over all bandit algorithms playing (I_1, I_2, \dots, I_T) , supremum is over all i.i.d. bernoulli reward distributions and expectation is over randomness of both rewards and algorithm.

We will prove the theorem in the next lecture.

Bibliography

- [1] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.