

Lecture 20: PAC(Probably Approximately Correct) learning II

*Instructors: Susan Murphy and Ambuj Tewari**Scribe: Aniket Deshmukh*

1 Review

In last lecture we talked about:

- PAC Learning: we want our solution/estimates to be approximately correct with high probability at the end of learning.
- Idea of an algorithm about PAC learning for Multi-Armed Bandits
- Median Elimination (ϵ, δ) algorithm [EDMM02]

In this lecture, we continue our discussion on median elimination (ϵ, δ) algorithm. Specifically, we talk about the lemma that we need, to prove that median elimination (ϵ, δ) is (ϵ, δ) -PAC.

2 Important Lemma

A key lemma is that, with high probability, the drop in the optimality of the best arm within the remaining set of arms is not severe in going from one phase to the next.

Lemma 1. *We have $\mathbb{P}(\max_{a \in S_l} \mu_a \leq \max_{a \in S_{l+1}} \mu_a + \epsilon_l) \geq 1 - \delta_l$*

Proof. w.l.o.g let us look at the first round. We want to show that

$$P(\mu^* \leq \max_{a \in s_2} \mu_a + \epsilon_1) \geq 1 - \delta_1$$

Let's look at the complement,

$$P(\mu^* > \max_{a \in s_2} \mu_a + \epsilon_1)$$

$$\begin{aligned}
\mu^* > \max_{a \in s_2} \mu_a + \epsilon_1 &\iff \forall a \in s_2, a \text{ is not } \epsilon_1 \text{ optimal} \\
&\implies a^* \notin s_2, \text{ that means } a^* \text{ was eliminated} \\
&\implies a \in s_2, \hat{\mu}_a^1 \geq \hat{\mu}_{a^*}^1 \\
&\implies \text{that means following set is large enough } (\geq \frac{K}{2}) \\
b_1 &= \{a \in \mathcal{A} : \mu_a < \mu^* - \epsilon_1 \text{ and } \hat{\mu}_a^1 \geq \hat{\mu}_{a^*}^1\}
\end{aligned}$$

Let bad event be the number of arms which are not ϵ -optimal but are empirically better than the best arm. Let $E_1 = \{\hat{\mu}_{a^*}^1 < \mu^* - \frac{\epsilon_1}{2}\}$

$$\begin{aligned}
P(\text{bad event}) &\leq P(b_1 \geq \frac{K}{2}) \\
&\leq P(E_1)P(b_1 \geq \frac{K}{2} \mid E_1) + P(E_1^c)P(b_1 \geq \frac{K}{2} \mid E_1^c) \\
&\leq P(b_1 \geq \frac{K}{2} \mid E_1^c) + P(E_1)
\end{aligned} \tag{1}$$

Now let's consider two terms in (1) separately. First let's bound the second term,

$$\begin{aligned}
P(E_1) &= \exp\left(-\left(\frac{\epsilon_1}{2}\right)^2 \left(\frac{1}{\epsilon_1/2}\right)^2 \log\left(\frac{3}{\delta_1}\right)\right) \quad (\text{Chernoff Hoeffding}) \\
&= \frac{\delta_1}{3}
\end{aligned} \tag{2}$$

Now for the first term,

$$\begin{aligned}
P(b_1 \geq \frac{K}{2} \mid E_1^c) &\leq \frac{\mathbb{E}[b_1 \mid E_1^c]}{K/2} \quad (\text{Markov inequality}) \\
&\leq \frac{K\delta_1/3}{K/2} = \frac{2\delta_1}{3} \quad (\text{We will prove this next})
\end{aligned} \tag{3}$$

Now all we need to show is $\mathbb{E}[b_1 \mid E_1^c] \leq \frac{K\delta_1}{3}$

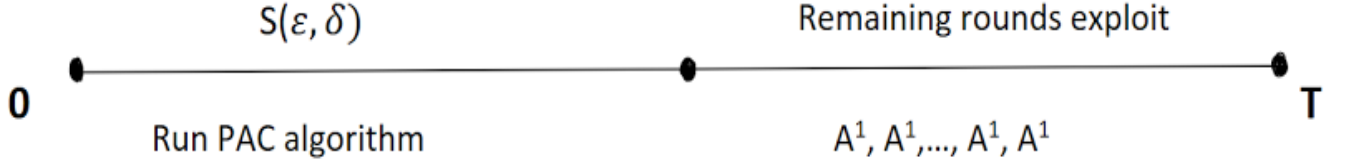
$$\begin{aligned}
\mathbb{E}[b_1 \mid E_1^c] &= \mathbb{E}\left[\sum_{a \in A, a \text{ isn't } \epsilon_1 \text{ optimal}} \mathbb{1}_{[\hat{\mu}_a^1 \geq \hat{\mu}_{a^*}^1 \mid E_1^c]}\right] \\
&= \sum_{a \in A, a \text{ isn't } \epsilon_1 \text{ optimal}} P(\hat{\mu}_a^1 \geq \hat{\mu}_{a^*}^1 \mid E_1^c)
\end{aligned} \tag{4}$$

Now consider an action which isn't ϵ_1 optimal.

$$\begin{aligned}
P(\hat{\mu}_a^1 \geq \hat{\mu}_{a^*}^1 \mid E_1^c) &= P(\hat{\mu}_a^1 \geq \hat{\mu}_{a^*}^1 \mid \hat{\mu}_{a^*}^1 \geq \mu^* - \frac{\epsilon_1}{2}) \\
&\leq P(\hat{\mu}_a^1 \geq \mu^* - \frac{\epsilon_1}{2} \mid \hat{\mu}_{a^*}^1 \geq \mu^* - \frac{\epsilon_1}{2}) \\
&\leq P(\hat{\mu}_a^1 \geq \mu_a + \frac{\epsilon_1}{2} \mid \hat{\mu}_{a^*}^1 \geq \mu^* - \frac{\epsilon_1}{2}) \\
&= P(\hat{\mu}_a^1 \geq \mu_a + \frac{\epsilon_1}{2}) \\
&\leq \exp\left(-\left(\frac{\epsilon_1}{2}\right)^2 \left(\frac{1}{\epsilon_1/2}\right)^2 \log\left(\frac{3}{\delta_1}\right)\right) \\
&= \frac{\delta_1}{3}
\end{aligned} \tag{5}$$

Using (5) and (4) we can conclude (3). Finally from (2) and (3) we can complete the proof of lemma. \square

Figure 1: PAC MAB algorithm



3 Relation between PAC setting and regret bound

Suppose we are given a (ϵ, δ) PAC algorithm with sample complexity $S(\epsilon, \delta)$. We could run this algorithm until $S(\epsilon, \delta)$ and then start exploiting the best arm (let's call it A^1) until remaining time $T - S(\epsilon, \delta)$.

Let R_T be the total regret until time T.

$$R_T = S(\epsilon, \delta) + \epsilon(T - S(\epsilon, \delta))(1 - \delta) + \delta(T - S(\epsilon, \delta)) \quad (6)$$

If $S(\epsilon, \delta) = \frac{k \log(\frac{1}{\delta})}{\epsilon^2}$ and $\delta = \frac{1}{T}$, then

$$\begin{aligned} R_T &= S(\epsilon, \delta) + \epsilon(T - S(\epsilon, \delta))(1 - \delta) + \delta(T - S(\epsilon, \delta)) \\ &\leq \frac{k \log(\frac{1}{\delta})}{\epsilon^2} + (1 - \delta)\epsilon T + \delta T \\ &= \frac{k \log(T)}{\epsilon^2} + \epsilon T + 1 \end{aligned} \quad (7)$$

by tuning ϵ we get total regret of the order $O(T^{2/3})$.

References

- [EDMM02] Eyal Even-Dar, Shie Mannor, and Yishay Mansour. Pac bounds for multi-armed bandit and markov decision processes. In *International Conference on Computational Learning Theory*, pages 255–270. Springer, 2002.