

Recall $\mu^* > \mu(v)$

$$\text{dinf}(v, \mu^*, \mathcal{E}) = \inf \{ D(v_i, v_i') : v_i' \in \mathcal{E}, \mu(v_i') > \mu^* \} \quad ①$$

π is consistent if $R_n(\pi, v) = o(n^a) \forall a > 0, v \in \mathcal{E}^k$

Thm: Given consistent policy π on \mathcal{E}^k , given v in \mathcal{E}^k

$$\liminf_{n \rightarrow \infty} \frac{R_n(\pi, v)}{\log n} \geq \sum_{i: \Delta_i > 0} \frac{\Delta_i}{\text{dinf}(v_i, \mu^*, \mathcal{E})} \quad ②$$

where μ^* is the mean of the optimal arm

① Here v is a scalar (distribution of a single arm)

② Here v is a vector (distributions of all arms, so an instance)

• Rich families have consistent policies, although they may not have logarithmic regret. Something like $(\log)^{(1+\epsilon)}$ if possible.

Proof: Suffices to show that for suboptimal i ,

$$\liminf_{n \rightarrow \infty} \frac{\mathbb{E}_v[T_i(n)]}{\log n} \geq \frac{1}{d_i}$$

where $d_i = \text{dinf}(v_i, \mu^*, \mathcal{E})$

$$\exists v_i' \in \mathcal{E} \text{ s.t. } \mu(v_i') > \mu^*, D(v_i, v_i') \leq d_i + \epsilon \in$$

The above follows from the definition of dinf

Define alternative instance $v' = (v_1, v_2, \dots, v_i', \dots, v_k)$ such that the i^{th} arm is changed and v_i' is as defined above.

We now define A s.t. it is a bad event in one instance & A^c is bad in the other

$$A = \{ T_i(n) > \frac{n}{2} \}$$

$$R_n(\pi, v) \geq \mathbb{P}_v(A) \cdot \frac{n \Delta_i}{2}$$

$$\text{In the alternative instance, } R_n(\pi, v') \geq \mathbb{P}_{v'}(A^c) \frac{n}{2} (\mu_i' - \mu^*)$$

$$\therefore R_n(\pi, v) + R_n(\pi, v') \geq c \cdot n [\mathbb{P}_v(A) + \mathbb{P}_{v'}(A^c)]$$

$$\text{where } c = \min \left(\frac{\Delta_i}{2}, \frac{\mu_i' - \mu^*}{2} \right)$$

Using the BH inequality,

$$R_n(\pi, v) + R_n(\pi, v') \geq \frac{cn}{2} \exp(-D(P_v, P_{v'}))$$

$$D(P_v, P_{v'}) = \sum_j \mathbb{E}_v[T_j(n)] D(v_j, v_j') = \mathbb{E}_v[T_i(n)] D(v_i, v_i')$$

$$D(v_i, v_i') \leq d_i + \epsilon$$

$$\therefore R_n(\pi, v) + R_n(\pi, v') \geq c'n \exp(-(d_i + \epsilon) \mathbb{E}_v[T_i(n)])$$

where $c' = c/2$

$$\frac{\log[R_n(\pi, v) + R_n(\pi, v')]}{\log n} \geq \frac{\log c'}{\log n} + \frac{1}{\log n} - \frac{(d_i + \epsilon) \mathbb{E}_v[T_i(n)]}{\log n} \quad (1)$$

Now we know that $R_n(\pi, v) = o(n^a) \forall a > 0$

Hence $\log[R_n(\pi, v)] = o(\log n^a) \xrightarrow{a \rightarrow 0} 0$ } EXERCISE

This is not obvious, because, for example $f = o(n) \not\Rightarrow \log f = o(\log n)$ when $f = \sqrt{n}$

Taking \liminf of both sides in (1),

$$\liminf_{n \rightarrow \infty} \frac{\mathbb{E}_v[T_i(n)]}{\log n} \geq \frac{1}{d_i + \epsilon}$$

Since this holds $\forall \epsilon > 0$, $\liminf_{n \rightarrow \infty} \frac{\mathbb{E}_v[T_i(n)]}{\log n} \geq \frac{1}{d_i}$

Note: d_i tries to capture how easy it is to confuse this arm with the optimal arm. If d_i is low, the possibility of confusion is high.

Note

"LEARNABILITY": $d_i > 0$

We can have $d_i = 0$ if

- \mathcal{E} = set of all bounded distributions (don't know $[a, b]$ a priori)
- \mathcal{E} = set of all subGaussian variable (σ is not known)

Here, we cannot estimate σ^2 along the way and still achieve logarithmic regret

This concept is known as "Statistical Robustness of the Environment"

BEST ARM IDENTIFICATION (PURE EXPLORATION)

In this paradigm, the learning is NOT minimization of regret, BUT identifying the best arm.

- 1) Fixed budget: The number of pulls is fixed. The algorithm outputs the arm it thinks is the best.
- 2) Fixed confidence: First fix the accuracy. The algorithm can make the mistake w.p. $\leq \delta$. Here we evaluate the algorithm based on how long it took for identification

These two are duals but not interconvertible

It turns out that the theory is easier for the fixed confidence setting

FIXED BUDGET SETTING

Given a budget of T pulls, where T is the horizon.

The algorithm gives an output $\hat{i} \in \{1, 2, \dots, k\}$. Single arm i , not a distribution.

We want to minimize $P(\hat{i} \neq i^*)$.

Model: k arms

② All 1-sub Gaussian

③ Horizon T

④ There is a unique best arm i^*

Note that assumption ④ is quite an important assumption in fixed ^{confidence} ~~budget~~ settings. It is just written for convenience here.

Note: We can also design a metric which penalizes instead of $P(\hat{i} \neq i^*)$ by the amount of suboptimality of the selected arm.

UNIFORM EXPLORATION ALGORITHM

• Pull each arm $\lfloor \frac{T}{k} \rfloor$ times.

• Output $\arg\max_i \hat{\mu}_i$

We shall try to bound the probability of an error of UEA.

Assume that arm 1 is the optimal arm.

$$P(\hat{i} \neq 1) = \sum_{j \neq 1} P(\hat{i} = j)$$

$$\therefore P(\hat{i} \neq 1) \leq \sum_{j \neq 1} P(\hat{\mu}_j \geq \hat{\mu}_1)$$

The inequality is because the condition is necessary but not sufficient.

$$\therefore P(\hat{i} \neq 1) \leq \sum_{j \neq 1} P(|\hat{\mu}_j - \mu_j| \geq \frac{\Delta_j}{2}) + P(|\hat{\mu}_1 - \mu_1| \geq \frac{\Delta_1}{2})$$

$$\leq 2(k-1) \exp\left(-\frac{\lfloor T/k \rfloor \Delta_2^2}{8}\right)$$

WLOG $\Delta_2 = \min(\Delta_i : \Delta_i \neq 0)$. Δ_2 is the term with the smallest decay rate

The probability of error decays exponentially with respect to the horizon.

We hence care about the decay rate, which should be as large as possible.

Here, the decay rate is $\frac{\Delta_2^2}{8k} = \frac{(\min \Delta_i)^2}{8k}$

Let us call the decay rate as α .

$$\therefore \alpha(U \# EA) = \frac{\Delta_2^2}{8K}$$

The simplest algorithm can be considered as a default baseline & now we want algorithm which have higher baselines.

Note: It turns out that superexponential decay is not possible. For example, something like $\exp(-\alpha T^2)$ is not possible.