

1 Lower bound on regret for K-armed bernoulli bandits

For any horizon T and number of arms K , there exists a strategy for choosing rewards such that the expected regret for any algorithm is $\Omega(\sqrt{TK})$. This strategy is oblivious to the algorithm and assigns rewards at random according to some distribution. Thus the regret bound applies to stochastic and adversarial bandits.

Theorem 1. *There exists a distribution over rewards such that:*

$$R(T) \geq \frac{1}{20} \min \{\sqrt{KT}, T\} \quad (1)$$

Proof. Consider a set of environments indexed by i . In environment i action i is the 'good' action. All other arms are slightly worse.

- In environment i , $r_{i,t} \sim \text{Bernoulli}(\frac{1}{2} + \epsilon)$ and $r_{j,t} \sim \text{Bernoulli}(\frac{1}{2}) \quad \forall j \neq i$
- $P_i(\cdot)$ is the probability with respect to environment i
- $P_{unif}(\cdot)$ is the probability with respect to an environment in which the expected reward for **all** arms is $\frac{1}{2}$
- $P_*(\cdot)$ is the probability with respect to an environment sampled uniformly at random from $\{1 \dots K\}$
- r_t is the reward received at time t
- $\mathbf{r}^t = \langle r_1, \dots, r_t \rangle$ is the history of rewards received upto time t
- A is the algorithm, maps $\mathbf{r}^{t-1} \rightarrow i_t$
- N_i is a random variable representing the number of times action i is selected by the algorithm.
- $G_A = \sum_{t=1}^T r_t$ is the total reward for the algorithm
- G_{max} is the total reward for playing the optimal arm in every timestep.

Lemma 2. *For any arm i ,*

$$\mathbb{E}_i[N_i] - \mathbb{E}_{unif}[N_i] \leq \frac{T}{2} \sqrt{-\ln(1 - 4\epsilon^2) \mathbb{E}_{unif}[N_i]} \quad (2)$$

This says that the number of times we expect to play arm i in the environment in which it is optimal is not too much greater than the number of times we expect to play it if all arms are equal.

Proof.

$$\mathbb{E}_i[N_i] - \mathbb{E}_{unif}[N_i] = \sum_{\mathbf{r} \in \{0,1\}^T} N_i(\mathbf{r}) (P_i(\mathbf{r}) - P_{unif}(\mathbf{r})) \quad \leftarrow \text{definition of expectation} \quad (3)$$

$$\leq \sum_{\mathbf{r}: P_i(\mathbf{r}) \geq P_{unif}(\mathbf{r})} N_i(\mathbf{r}) (P_i(\mathbf{r}) - P_{unif}(\mathbf{r})) \quad \leftarrow \text{dropped only -ive terms} \quad (4)$$

$$\leq T \sum_{\mathbf{r}: P_i(\mathbf{r}) \geq P_{unif}(\mathbf{r})} (P_i(\mathbf{r}) - P_{unif}(\mathbf{r})) \quad \leftarrow N_i \leq T \quad \forall \mathbf{r} \quad (5)$$

$$= \frac{T}{2} \|P_i(\mathbf{r}) - P_{unif}(\mathbf{r})\|_1 \quad \leftarrow \text{see Thomas\&Cover 11.137} \quad (6)$$

$$\leq \frac{T}{2} \sqrt{2 \ln(2) KL(P_{unif}(\mathbf{r}) \| P_i(\mathbf{r}))} \quad \leftarrow \text{see Thomas\&Cover 11.138} \quad (7)$$

$$= \frac{T}{2} \sqrt{-\ln(2) \lg(1 - 4\epsilon^2) \mathbb{E}_{unif}[N_i]} \quad \leftarrow \text{see section 1.1} \quad (8)$$

$$= \frac{T}{2} \sqrt{-\ln(1 - 4\epsilon^2) \mathbb{E}_{unif}[N_i]} \quad \leftarrow \text{change of base} \quad (9)$$

□

Theorem 3. For any algorithm A , if the distribution over rewards is selected uniformly at random from environments $\{1 \dots K\}$:

$$E_*[G_{max} - G_A] \geq \epsilon \left(T - \frac{T}{K} - \frac{T}{2} \sqrt{-\frac{T}{K} \ln(1 - 4\epsilon^2)} \right) \quad (10)$$

Proof.

$$\mathbb{E}_i[r_t] = \left(\frac{1}{2} + \epsilon \right) P_i(i_t = i) + \frac{1}{2} (1 - P_i(i_t = i)) \quad (11)$$

$$= \frac{1}{2} + \epsilon P_i(i_t = i) \quad (12)$$

$$\implies \mathbb{E}_i[G_A] = \sum_{t=1}^T \mathbb{E}_i[r_t] = \frac{T}{2} + \epsilon \mathbb{E}_i[N_i] \quad (13)$$

This gives us the expected gain given action i is the good action in terms of the number of times the algorithm A selects action i . The expected gain over all the environments i is:

$$\mathbb{E}_*[G_A] = \frac{1}{K} \sum_{i=1}^K \mathbb{E}_i[G_A] \quad (14)$$

$$= \frac{T}{2} + \frac{\epsilon}{K} \sum_{i=1}^K \mathbb{E}_i[N_i] \quad (15)$$

$$\mathbb{E}_*[G_{max}] = \left(\frac{1}{2} + \epsilon \right) T \quad (16)$$

From lemma 2

$$\sum_{i=1}^K \mathbb{E}_i[N_i] \leq \sum_{i=1}^K \left(\mathbb{E}_{unif}[N_i] + \frac{T}{2} \sqrt{-\ln(1 - 4\epsilon^2) \mathbb{E}_{unif}[N_i]} \right) \quad (17)$$

Now $\sum_{i=1}^K \mathbb{E}_{unif}[N_i] = T$ (because $\sum_{i=1}^K N_i = T$? but doesn't that imply $\sum_{i=1}^K \mathbb{E}_i[N_i] = T$?)

$$\implies \sum_{i=1}^K \mathbb{E}_i[N_i] \leq T + \frac{T}{2} \sqrt{-\ln(1 - 4\epsilon^2) K T} \quad \leftarrow \text{via Jensen's Inequality} \quad (18)$$

$$\implies \mathbb{E}_*[G_A] \leq \frac{T}{2} + \epsilon \left(\frac{T}{K} + \frac{T}{2} \sqrt{-\frac{T}{K} \ln(1 - 4\epsilon^2)} \right) \quad (19)$$

$$\implies \mathbb{E}_*[G_{max} - G_A] \geq \epsilon \left(T - \frac{T}{K} - \frac{T}{2} \sqrt{-\frac{T}{K} \ln(1 - 4\epsilon^2)} \right) \quad (20)$$

□

□

1.1 Calculation of KL divergence

$$KL(p(x_1, \dots, x_n) || q(x_1, \dots, x_n)) = \sum_{i=1}^N KL(p(x_i | \mathbf{x}^{i-1}) || q(x_i | \mathbf{x}^{i-1})) \quad \leftarrow \text{chain rule for KL divergence} \quad (21)$$

$$\text{where } KL(p(x_i | \mathbf{x}^{i-1}) || q(x_i | \mathbf{x}^{i-1})) = \sum_{\mathbf{x}^{i-1}} p(\mathbf{x}^{i-1}) \sum_{x_i} p(x_i | \mathbf{x}^{i-1}) \lg \left(\frac{p(x_i | \mathbf{x}^{i-1})}{q(x_i | \mathbf{x}^{i-1})} \right) \quad (22)$$

So

$$KL(P_{unif} || P_i) = \sum_{t=1}^T \left(\underbrace{\sum_{\mathbf{r}^{t-1}} P_{unif}(\mathbf{r}^{t-1})}_{\text{expectation over history}} \underbrace{\sum_{r_t \in \{0,1\}} P_{unif}(r_t | \mathbf{r}^{t-1}) \lg \left(\frac{P_{unif}(r_t | \mathbf{r}^{t-1})}{P_i(r_t | \mathbf{r}^{t-1})} \right)}_{\text{KL divergence at time t}} \right) \quad (23)$$

Now

$$P_{unif}(r_t | \mathbf{r}^{t-1}) = \frac{1}{2} \quad \forall r_t, \mathbf{r}^{t-1} \quad (24)$$

$$P_i(r_t | \mathbf{r}^{t-1}) = \begin{cases} (\frac{1}{2} + \epsilon)^{r_t} + (\frac{1}{2} - \epsilon)^{1-r_t} & \text{if } A(\mathbf{r}^{t-1}) = i \\ \frac{1}{2} & \text{otherwise} \end{cases} \quad (25)$$

Let B be the set of histories that lead the algorithm to select the good arm, $B = \{\mathbf{r}^{t-1} : A(\mathbf{r}^{t-1}) = i\}$

$$KL(P_{unif} || P_i) = \sum_{t=1}^T \left(\sum_B P_{unif}(\mathbf{r}^{t-1}) KL\left(\frac{1}{2} || \frac{1}{2} + \epsilon\right) + \sum_{B^c} P_{unif}(\mathbf{r}^{t-1}) KL\left(\frac{1}{2} || \frac{1}{2}\right) \right) \quad (26)$$

$$= KL\left(\frac{1}{2} || \frac{1}{2} + \epsilon\right) \sum_{t=1}^T \left(\sum_B P_{unif}(\mathbf{r}^{t-1}) \right) \quad (27)$$

$$= KL\left(\frac{1}{2} || \frac{1}{2} + \epsilon\right) \sum_{t=1}^T (P_{unif}(i_t = i)) \quad (28)$$

$$= KL\left(\frac{1}{2} || \frac{1}{2} + \epsilon\right) \mathbb{E}_{unif}[N_i] \quad (29)$$

$$= -\frac{1}{2} \lg(1 - 4\epsilon^2) \mathbb{E}_{unif}[N_i] \quad (30)$$

1.2 Jenson's inequality

Jenson's inequality states that for a concave function ϕ :

$$\frac{\sum_{i=1}^K \phi(x_i)}{K} \leq \phi\left(\frac{\sum_{i=1}^K x_i}{K}\right) \quad (31)$$

$$\implies \sum_{i=1}^K \sqrt{\mathbb{E}_{unif}[N_i]} \leq \sqrt{KT} \quad (32)$$

$$\implies \sum_{i=1}^K \mathbb{E}_i[N_i] \leq T + \frac{T}{2} \sqrt{-\ln(1 - 4\epsilon^2)KT} \quad (33)$$