

# Multi-armed bandits

Stochastic: Pulling arm  $i$  results in R.V.  $X_i$  w/  $E[X_i] = \mu_i$  <sup>reward</sup>  $\forall$  time

$$\Delta_i = \mu_{i^*} - \mu_i$$

$$\max_i E \left[ \sum_{t=1}^T \mu_i - \mu_{I_t} \right] \leq \sqrt{nT \log(nT)} \wedge \sum_{i \neq i^*} \Delta_i^{-1} \log(nT)$$

Adversarial setting: Adversary chooses  $l_t \in [-1, 1]^n$   $\forall t$  prior to start of game.

Pulling arm  $i$  results in loss  $l_{t,i}$

$$\max_i E \left[ \sum_{t=1}^T l_{t,I_t} - l_{t,i} \right] \leq \sqrt{nT \log(nT)}$$

Thm 1 For every  $T \geq n$   $\exists$  set of bandit instances w/  $P_\theta = \mathcal{N}(\theta, I)$

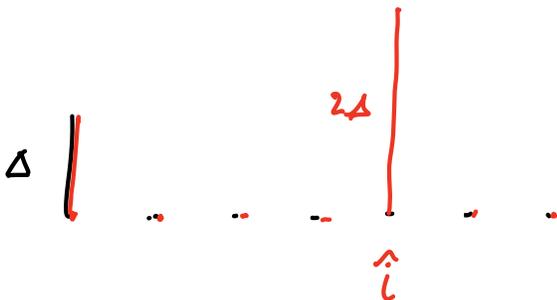
s.t.  $\sup_{P_\theta} E_\theta [\text{Regret}] \geq \sqrt{(n-1)T/256}$

$$\theta = (\Delta, 0, \dots, 0) \in [0, 1]^n$$

$$\theta' = (\Delta, 0, \dots, 2\Delta, \dots, 0)$$

$$\exists \hat{i} \in \{2, \dots, n\} \text{ s.t. } E_\theta [T_{\hat{i}}] \leq \frac{T}{n-1}$$

$\Rightarrow$  intuitively, we can't tell whether  $\theta_{\hat{i}} = 0$  or  $\theta_{\hat{i}} \approx \sqrt{\frac{n-1}{T}} \approx \Delta$



$$\text{Regret}(\theta') = \sum_i \Delta_i E_{\theta'} [T_i] \geq \Delta_i E_{\theta'} [T_i] \geq \Delta P_{\theta'}(T_i \geq T/2) T/2$$

$$\text{Regret}(\theta) = \sum_{i \neq 1} \Delta_i \mathbb{E}_\theta [T_i] = \Delta \sum_{i \neq 1} \mathbb{E}_\theta [T_i]$$

$$\geq \Delta \mathbb{P}_\theta \left( \sum_{i \neq 1} T_i \geq T/2 \right) T/2$$

$$= \Delta \left( 1 - \mathbb{P}_\theta (T_1 \geq T/2) \right) T/2$$

$$\max_{\mu \in (\theta, \theta')} \text{Regret}(\mu) \geq \frac{1}{2} \left( \text{Regret}(\theta) + \text{Regret}(\theta') \right)$$

$$\geq \frac{T\Delta}{4} \left( 1 + \mathbb{P}_{\theta'} (T_1 \geq T/2) - \mathbb{P}_\theta (T_1 \geq T/2) \right)$$

$$\geq \frac{\Delta T}{4} \left( 1 - \sup_A |\mathbb{P}_{\theta'}(A) - \mathbb{P}_\theta(A)| \right)$$

(Pinsker's ineq.)  $\geq \frac{\Delta T}{4} \left( 1 - \sqrt{\text{KL}(\mathbb{P}_\theta \parallel \mathbb{P}_{\theta'}) / 2} \right)$

$$= \frac{\Delta T}{4} \left( 1 - \sqrt{\sum_i (\theta_i - \theta'_i)^2 \mathbb{E}_\theta [T_i] / 4} \right)$$

$$= \frac{\Delta T}{4} \left( 1 - \sqrt{\Delta^2 \mathbb{E}_\theta [T_i]} \right)$$

$$\geq \frac{\Delta T}{4} \left( 1 - \sqrt{\Delta^2 T / (n-1)} \right)$$

$$\Delta = \sqrt{\frac{n-1}{4T}}$$

$\mathbb{P}_{\theta'}, \mathbb{P}_\theta$  defined over  $\mathcal{I}_d$



$$\sup_A |\mathbb{P}_{\theta'}(A) - \mathbb{P}_\theta(A)| = \int_{\mathcal{X}} \left| \frac{d\mathbb{P}_\theta(x)}{dx} - \frac{d\mathbb{P}_{\theta'}(x)}{dx} \right| dx$$

$$\leq \int_x \frac{dP_\theta}{dx} \log \left( \frac{\frac{dP_\theta(x)}{dx}}{\frac{dP_{\theta'}(x)}{dx}} \right) dx$$

$$= \text{KL}(P_\theta \parallel P_{\theta'})$$

□

Consider a game where at each time  $t$

row player chooses  $I_t \in [m]$ , column player chooses  $J_t \in [n]$

and row player receives reward  $e_{I_t}^T A e_{J_t} + \gamma_t$

where  $A$  is unknown and  $\mathbb{E}[\gamma_t] = 0$ ,  $|\gamma_t| \leq 1$ .

Row-player's objective maximize  $\mathbb{E} \left[ \sum_t e_{I_t}^T A e_{J_t} \right]$

$$l_t = -A e_{J_t}$$

$V_A = \max_{x \in \Delta_m} \min_{y \in \Delta_n} x^T A y$  is the value of the game.

$(x, y) \in \Delta_m \times \Delta_n$  is a  $\epsilon$ -Nash equilibrium if they achieve the value of the game. Equivalently:

$$(x - x')^T A y \geq -\epsilon \quad \forall x' \quad \text{and} \quad x^T A (y' - y) \geq -\epsilon \quad \forall y'$$

$$\text{Nash-Regret} = \mathbb{E} \left[ \sum_{t=1}^T (V_* - e_{i_t}^T A e_{j_t}) \right] \quad \begin{array}{l} I_t \sim x_t \\ J_t \sim y_t \end{array}$$

$$= \mathbb{E} \left[ \sum_{t=1}^T x_t^T A y_0 - x_t^T A y_t \right]$$

$$\leq \mathbb{E} \left[ \sum_{t=1}^T x_t^T A y_t - x_t^T A y_t \right]$$

$$= \mathbb{E} \left[ \sum_{t=1}^T (x_t - x_t^*)^T A y_t \right] = \text{External regret.}$$

Claim]  $\exists$  A matrix s.t. for any row-player that achieves  $o(T)$  external regret,  $\exists$  adversary s.t.

$$\text{external regret} \geq c\sqrt{mT} \text{ and } \mathbb{E} \left[ \sum_{t=1}^T x_t^T A y_t - x_t^* A y_t \right] \geq c'T$$

Let row-player and column player each play an independent copy of EXP-3. and let  $\hat{x}_T = \frac{1}{T} \sum_{t=1}^T e_{i_t}$ ,  $\hat{y}_T = \frac{1}{T} \sum_{t=1}^T e_{j_t}$

$$(x - \hat{x})^T A \hat{y} \leq c_1 \sqrt{\frac{m \log(T)}{T}} \quad \hat{x}^T A (y - \hat{y}) \leq c \sqrt{\frac{n \log(T)}{T}}$$

$\Rightarrow$  if  $T \geq c'(m+n) \bar{\epsilon}^{-2} \log(1/\epsilon)$  then  $(\hat{x}, \hat{y})$  is an  $\epsilon$ -Nash Eq.

$$\sum_t e_{J_t}^T A e_{J_t} = \underbrace{\sum_t (e_{J_t} - x)^T A e_{J_t}}_{\geq -c\sqrt{nT \log T}} + \underbrace{\sum_t x^T A e_{J_t}}_{T \cdot x^T A \hat{y}_T}$$

$$\geq -c\sqrt{nT \log T} + T \cdot x^T A \hat{y}_T$$

$$R \begin{bmatrix} R & P & S \\ 0 & -1 & 1 \\ P & 1 & 0 & -1 \\ S & -1 & 1 & 0 \end{bmatrix}$$

$$\sum_t e_{J_t}^T A e_{J_t} = \sum_t e_{J_t}^T A (e_{J_t} - y) + \sum_t e_{J_t}^T A y$$

$$\leq c\sqrt{nT \log T} + T \cdot \hat{x}_T^T A y$$



$$x^T A \hat{y}_T - \hat{x}_T^T A y \leq 2c\sqrt{\frac{n \log T}{T}}$$

$$\Rightarrow \max_{x, y} (x - \hat{x}_T)^T A \hat{y}_T + \hat{x}_T^T A (\hat{y}_T - y) \leq 2c\sqrt{\frac{n \log T}{T}}$$

$$\Rightarrow \varepsilon = 2c\sqrt{\frac{n \log T}{T}}$$

$$(x - \hat{x}_T)^T A \hat{y}_T \leq \varepsilon$$

$$\hat{x}_T^T A (\hat{y}_T - y) \leq \varepsilon$$

$$\min_y \max_x x^T A y - c \sqrt{\frac{n \log T}{T}}$$

$$\leq \max_x x^T A \hat{y}_T - c \sqrt{\frac{n \log T}{T}}$$

$$\leq \min_y \hat{x}_T^T A y + c \sqrt{\frac{n \log T}{T}}$$

$$\leq \max_x \min_y x^T A y + c \sqrt{\frac{n \log T}{T}}$$

$$\min_y \max_x x^T A y \leq \max_x \min_y x^T A y$$

---

Sion's minimax theorem