

EXP3(γ): Exponential Weights for Exploration Exploitation

Input: Time horizon T , n arms, $\eta > 0$, $\gamma \in [0, 1]$, $\lambda \in \Delta_n$.

Initialize: Player sets $p_1 = (1/n, \dots, 1/n) \in \Delta_n$. Adversary chooses $\{\ell_t\}_{t=1}^T \subset [-1, 1]^n$.

for: $t = 1, \dots, T$

Player draws $I_t \sim q_t := (1 - \gamma)p_t + \gamma\lambda$ and suffers (and observes) loss $\ell(I_t, \ell_t) = \ell_{t,I_t}$

Player computes $\hat{\ell}_{t,i} = \frac{\mathbf{1}\{I_t=i\}}{q_{t,i}} \ell_{t,i}$

Update iterates:

$$w_{t+1,i} = w_{t,i} \exp(-\eta \hat{\ell}_{t,i}) \quad p_{t+1,i} = w_{t+1,i} / \sum_{j=1}^n w_{t+1,j}.$$

Theorem 17. Fix any sequence $\ell_t \in [-1, 1]$ for all t and assume $\eta, \gamma \geq 0$ satisfy $-\eta \hat{\ell}_{t,i} \leq 1$ for all i, t . Then we have

$$\gamma = 3d$$

$$\max_{i=1, \dots, n} \mathbb{E} \left[\sum_{t=1}^T \ell_{t,I_t} - \ell_{t,i} \right] \leq 2\gamma T + \frac{\log(n)}{\eta} + (1 - \gamma)\eta \mathbb{E} \left[\sum_{t=1}^T \sum_{j=1}^n p_{t,j} \hat{\ell}_{t,j}^2 \right].$$

$$= 2\gamma dT + \frac{\log(n)}{\gamma} + \gamma dT \frac{dT}{1-\gamma} = 3\gamma dT + \frac{\log(n)}{\gamma} \leq \sqrt{12dT \log(n)}$$

Linear Bandits - Adversarial Setting

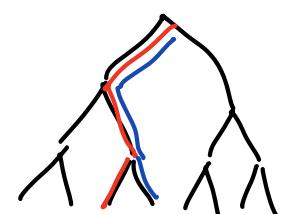
Input $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$

Adversary chooses $\{\theta_t\}_{t=1}^T \subset \mathbb{R}^d$: $\max_{x \in X} |\langle x, \theta_t \rangle| \leq 1 \forall t$

for $t = 1, 2, \dots, T$

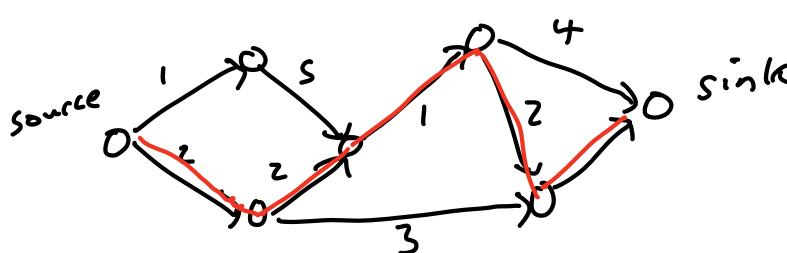
Player chooses $x_t \in X$

Adversary reveals $\langle \theta_t, x_t \rangle$ (not θ_t !)



$$\text{Regret} = \max_{x \in X} \sum_{t=1}^T \langle \theta_t, x_t - x \rangle$$

x	0	x
x	0	



$$x \in \{0, 1\}^{IE}$$

For $x_i \in \mathcal{X}$ define $\ell_{t,i} = \langle x_i, \theta_t \rangle$

$$\hat{\ell}_{t,i} = x_i^\top A(g_t)^{-1} x_{I_t} \ell_{t,I_t}, \quad A(\lambda) = \sum_{i=1}^n \lambda_i x_i x_i^\top$$

$$\begin{aligned} \mathbb{E}[A(g_t)^{-1} x_{I_t} \ell_{t,I_t}] &= \sum_{j=1}^n g_{t,j} A(g_t)^{-1} x_j \ell_{t,j} \\ &= \sum_{j=1}^n g_{t,j} A(g_t)^{-1} x_j x_j^\top \theta_t \\ &= I \cdot \theta_t = \theta_t \end{aligned}$$

$$\Rightarrow \mathbb{E}[\hat{\ell}_{t,i}] = \ell_{t,i} \quad \forall i$$

Aside $x_t \sim \lambda$, $y_t = \langle x_t, \theta_* \rangle + \underline{\eta}_t \quad t=1, \dots, T$

$$\hat{\theta}_{LS} = \left(\frac{1}{T} \sum_{t=1}^T x_t x_t^\top \right)^{-1} \frac{1}{T} \sum_{t=1}^T x_t y_t$$

$$\mathbb{E}[\hat{\theta}_{LS}] = \theta_*, \quad \text{Cov}(\hat{\theta}_{LS}) = \left(\sum_{t=1}^T x_t x_t^\top \right)^{-1}$$

$$\hat{\theta}_{IPs} = \left(\sum_{x \in \mathcal{X}} \lambda_x x x^\top \right)^{-1} \frac{1}{T} \sum_{t=1}^T x_t y_t$$

$$\mathbb{E}[\hat{\theta}_{IPs}] = \theta_*, \quad \text{Cov}(\hat{\theta}_{IPs}) = \left(T \sum_{x \in \mathcal{X}} \lambda_x x x^\top \right)^{-1}$$

$$\mathcal{X} = \{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \}$$

$$\lambda_1, \lambda_2 = 1 - \lambda,$$

$$T_1 = \sum_{t=1}^T \mathbb{I}\{\mathbf{I}_t = 1\}$$

$$\hat{\theta}_{LS} = \left[\frac{1}{T_1} \sum_{t=1}^{T_1} \mathbb{I}\{\mathbf{I}_t = 1\} y_t \right. \\ \left. \frac{1}{T-T_1} \sum_{t=T_1+1}^T \mathbb{I}\{\mathbf{I}_t = 2\} y_t \right]$$

$$\hat{\theta}_{IPS} = \left[\frac{1}{\lambda T} \sum_{t=1}^T \mathbb{I}\{\mathbf{I}_t = 1\} y_t \right. \\ \left. \frac{1}{(1-\lambda)T} \sum_{t=T_1+1}^T \mathbb{I}\{\mathbf{I}_t = 1\} y_t \right]$$

$$\tilde{l}_{t,i} = x_i^\top A(g_t)^{-1} x_{I_t} l_{t,I_t}$$

$$- \mathbb{E} \tilde{l}_{t,i} \leq \mathbb{E} |x_i^\top A(g_t)^{-1} x_{I_t}| \quad (|l_{t,i}| \leq 1)$$

$$= \mathbb{E} |\langle A(g_t)^{-1/2} x_i, A(g_t)^{-1/2} x_{I_t} \rangle|$$

$$\leq \mathbb{E} \|x_i\|_{A(g_t)^{-1}} \cdot \|x_{I_t}\|_{A(g_t)^{-1}}$$

$$\leq \max_{j=1,\dots,n} \mathbb{E} \cdot \|x_j\|_{A(g_t)^{-1}}^2$$

$$\leq \max_{j=1,\dots,n} \frac{3}{8} \|x_j\|_{A(\lambda)^{-1}}^2 \quad \text{if } \lambda \text{ G-optimal.}$$

$$= \frac{3}{8} \cdot d$$

Suggests picking $\gamma = 2d$ s.t. $-3\hat{L}_{t,i} \leq 1$ t.c.

$$\mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} \hat{L}_{t,i}^2 \right] = \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} x_i^T A(g_t)^{-1} x_{I_t} x_{I_t}^T A(g_t)^{-1} x_i \hat{L}_{t,I_t}^2 \right]$$

$$\leq \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} x_i^T A(g_t)^{-1} x_{I_t} x_{I_t}^T A(g_t)^{-1} x_i \right]$$

$$= \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} x_i^T A(g_t)^{-1} x_i \right]$$

$$\leq \frac{1}{1-\gamma} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} x_i^T A(P_t)^{-1} x_i \right]$$

$$= \frac{1}{1-\gamma} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} \text{Trace} \left(x_i^T A(P_t)^{-1} x_i \right) \right]$$

$$= \frac{1}{1-\gamma} \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^n P_{t,i} \text{Trace} \left(A(P_t)^{-1} x_i x_i^T \right) \right]$$

$$= \frac{dT}{1-\gamma}$$

Contextual Bandits

for $t=1, 2, \dots, T$

Nature reveals context $c_t \stackrel{iid}{\sim} \mathcal{V}$

Player chooses $I_t \in [n]$

and gets reward $y_t \in [-1, 1]$, s.t. $\mathbb{E}[y_t | c_t] = r(c_t, I_t)$

Suppose support(\mathcal{V}) = C and $|C| < \infty$.

Treat each context as indep. bandit.

$$\sum_{t=1}^T \max_i (r(c_t, i) - r(c_t, I_t))$$

$$= \sum_{c \in C} \max_i \sum_{t=1}^T \mathbb{I}\{c_t = c\} (r(c_t, i) - r(c_t, I_t))$$

$$\leq \sum_{c \in C} \sqrt{n T_c} \quad \leftarrow \text{If playing MAB w/ } \sqrt{nT} \text{ regret}$$

$$T_c := \sum_t \mathbb{I}\{c_t = c\}$$

$$\leq \sqrt{\left(\sum_{c \in C} (\sqrt{T_c})^2 \right) \cdot \left(\sum_{c \in C} (\sqrt{n})^2 \right)}$$

$$= \sqrt{n |C| T}$$

Alternatively, what if we ignore context?

$$\max_i \sum_{t=1}^T r(c_t, i) - r(c_t, I_t) \leq \sqrt{nT}$$