

Exp3 in the Linear Setting

Source: Bandit Algorithms, L. & S., Chap 27.

Adversarial Setting for Linear Bandits.

Setting: Actions $A \subseteq \mathbb{R}^d$, $t = 1, 2, \dots, n$

Player: Choose $A_t \in A$, and receives reward
$$Y_t = \langle A_t, y_t \rangle.$$

loss vector chosen by adversary; $y_t \in \mathbb{R}^d$

Adversary: As before, adversary aware of player's policy (but not action, as it is randomized using secret bits; see Set 8 - Exp3 Algorithm).

Assumptions:

① $\forall t = 1, 2, \dots, n$, $y_t \in \mathcal{A} = \{z \in \mathbb{R}^d : \sup_{a \in A} |\langle a, z \rangle| \leq 1\}$

② A spans \mathbb{R}^d

Regret: $R_n = E \left[\sum_{t=1}^n \gamma_t \right] - \min_{a \in \mathcal{A}} \sum_{t=1}^n \langle a, y_t \rangle$

Linear Exp3 Algorithm Idea: (supposing $|\mathcal{A}|$ finite)

(a) Based on $\gamma_t \in \mathbb{R}$, construct $\hat{\gamma}_t \in \mathbb{R}^d$, which is an unbiased estimate for $y_t \in \mathbb{R}^d$. (observed loss which is the projection of y_t along A_t)

(b) Use $\hat{\gamma}_t$ to exponentially reweight prob. distribution over $a \in \mathcal{A}$.

(c) Play (using secret randomness) a randomly chosen arm $a \in \mathcal{A}$ using this distribution.

Assumption: $|\mathcal{A}|$ is finite (for now).

Importance Sampling Estimator:

$$\hat{\gamma}_t = Q_t^{-1} A_t \gamma_t, \text{ where } Q_t = \sum_{a \in \mathcal{A}} \underbrace{P_t(a)}_{\substack{\text{prob. with} \\ \text{which arm } a \\ \text{played at time } t}} a a^T,$$

$$\hat{\gamma}_t(a) = \langle a, \hat{\gamma}_t \rangle.$$

$$P_t(a) = \underbrace{\gamma \pi(a)}_{\text{an exploration distribution that is used for variance control, } \pi \text{ chosen using Kiefer-Wolfowitz Thm}} + (1-\gamma) \frac{e^{-\eta \sum_{s=1}^{t-1} \hat{y}_s(a)}}{\sum_{b \in \mathcal{A}} e^{-\eta \sum_{s=1}^{t-1} \hat{y}_s(b)}}$$

an exploration distribution that is used

for variance control, π chosen using Kiefer-Wolfowitz Thm

Parameters: We will see $\gamma = g(\pi) \cdot \eta$, $\eta \sim \sqrt{\frac{\log(k)}{(2g(\pi)+d)n}}$
 \nearrow Kiefer-Wolfowitz
 $g(\pi) \leq d$

Algorithm: Input $A \subseteq \mathbb{R}^d$, η the learning rate, exploration dist. π , exploration parameter γ .

For each $t=1, 2, \dots, n$:

① Compute $P_t(a)$ as above, and play arm $A_t \sim P_t$.

② Update $\hat{y}_t = Q_t^{-1} A_t \gamma_t = Q_t^{-1} A_t \langle A_t, y_t \rangle$

with $\hat{y}_t(a) = \langle a, \hat{y}_t \rangle$

Theorem (27.1 in text): Setting as above. Then
 $\exists \gamma, \eta$ s.t. for any π ,

$$R_n \leq 2\sqrt{(2g(\pi) + d)n \log k},$$

where $k = |A|$, $g(\pi) = \max_{a \in A} \|a\|_{Q(\pi)^{-1}}^2$

Further, $\exists \pi$ s.t. $g(\pi) \leq d$ and

$$Q(\pi) = \sum_{a \in A} \pi(a) a a^T$$

$$R_n \leq 2\sqrt{3dn \ln(k)}$$

Proof: Similar in spirit to Exp3 proof. Assume
 η small enough s.t. $-1 \leq \eta \hat{Y}_t(a) \leq 1 \quad \forall a \in A$.

↳ explicit value later in proof.

Using Exp3 proof, we can get (analog of Step 1):

$$R_n \leq \frac{\ln(k)}{\eta} + 2\gamma n + \eta \sum_{t=1}^n \mathbb{E} \left[\sum_{a \in A} P_t(a) \hat{Y}_t^2(a) \right]$$

Bounding $\mathbb{E} \left[\underbrace{\sum_{a \in A} P_t(a) \hat{Y}_t^2(a)}_{\triangleq M_t} \right]$:

↓
 ①

↳ $\triangleq M_t$

$$\begin{aligned}\hat{\gamma}_t^2(a) &= (\langle a, \hat{\gamma}_t \rangle)^2 = (\langle a, Q_t^{-1} A_t \gamma_t \rangle)^2 \\ &= (a^T \underbrace{Q_t^{-1} A_t}_{\text{scalar}} \gamma_t)^2 = \gamma_t^2 (A_t^T Q_t^{-1} a a^T Q_t^{-1} A_t)\end{aligned}$$

$$\therefore M_t = \sum_{a \in A} P_t(a) \hat{\gamma}_t^2(a)$$

$$= \gamma_t^2 A_t^T Q_t^{-1} \underbrace{\left(\sum_{a \in A} P_t(a) a a^T \right)}_{Q_t} Q_t^{-1} A_t$$

$$= \underbrace{\gamma_t^2}_{\leq 1} A_t^T Q_t^{-1} A_t \leq A_t^T Q_t^{-1} A_t = \text{trace}(A_t^T Q_t^{-1} A_t)$$

$$= \text{trace}(A_t A_t^T Q_t^{-1})$$

$$\therefore E[M_t | A_1, A_2, \dots, A_{t-1}] \leq \text{trace} \left(\underbrace{\left(\sum_{a \in A} P_t(a) a a^T \right) Q_t^{-1}}_{Q_t} \right)$$

= d.

Recall: $|\eta \hat{\gamma}_t(a)| \leq 1$, and $|\gamma_t| \leq 1$. Thus,

$$|\eta \hat{\gamma}_t(a)| = |\eta a^T Q_t^{-1} A_t \gamma_t| \leq \eta |a^T Q_t^{-1} A_t| \quad \text{(*)}$$

Now, we have:

Recall

$$Q(\pi) = \sum_{a \in \mathcal{A}} \pi(a) a a^T$$

$$Q_t = \sum_{a \in \mathcal{A}} P_t(a) a a^T \geq \gamma Q(\pi)$$

$\gamma \pi(a) + (1-\gamma)(\cdot)$

$$\therefore Q_t^{-1} \preceq \frac{Q(\pi)^{-1}}{\gamma}$$

⊗

Thus, $|a^T Q_t^{-1} A_t| = |\langle a, Q_t^{-1} A_t \rangle| \leq \|a\|_{Q_t^{-1}} \|Q_t^{-1} A_t\|_{Q_t}$

$$= \|a\|_{Q_t^{-1}} (A_t^T Q_t^{-1} Q_t Q_t^{-1} A_t) = \|a\|_{Q_t^{-1}} \|A_t\|_{Q_t^{-1}}$$

$$\leq \max_{v \in \mathcal{A}} v^T Q_t^{-1} v \leq \frac{1}{\gamma} \max_{v \in \mathcal{A}} v^T Q(\pi)^{-1} v = \frac{g(\pi)}{\gamma}$$

recall set 16:

Kiefer-Wolfowitz Thm

$$\text{Thus, } |\eta \hat{Y}_t(a)| \leq \frac{\eta}{\gamma} g(\pi).$$

$$\text{Choose } \gamma = \eta g(\pi) \Rightarrow |\eta \hat{Y}_t(a)| \leq 1.$$

\therefore Substituting in ①:

$$R_n \leq \frac{\ln(k)}{\eta} + \underbrace{2\gamma n}_{\eta g(\pi)} + \underbrace{\eta \sum_{t=1}^n \mathbb{E} \left[\sum_{a \in A} P_t(a) \hat{Y}_t^2(a) \right]}_{\leq d}$$

$$\leq \frac{\ln(k)}{\eta} + \eta n (2g(\pi) + d)$$

$$\searrow = \sqrt{\frac{\ln(k)}{2g(\pi) + d}}$$

$$= 2 \sqrt{(2g(\pi) + d) n \ln(k)}$$

