

Exp 3 Algorithm

Source: Bandit Algorithms, L&S., Chap. 11

Adversarial Setting

Bandit system evolves over discrete time steps

Two actors: Adversary chooses arm rewards

Player chooses arm to play

k arms, $j = \{1, 2, \dots, k\}$, no stochastic assumption on arm rewards.

At each time $t = 1, 2, \dots, n$:

① Adversary chooses reward for each arm j :

$$x_{tj} \in [0, 1], j = 1, 2, \dots, n$$

equivalently chooses loss for each arm j : arbitrary /

$$y_{tj} \in [0, 1], j = 1, 2, \dots, n$$

$$y_{tj} = (1 - x_{tj})$$

no assumptions
on the
way these
are chosen

② Player chooses action $A_t \in \{1, 2, \dots, k\}$.

A_t chosen in a two step process:

① A policy $\pi = (\pi_1, \pi_2, \dots, \pi_n)$

is fixed before start of play, where

history of action, rewards

$$\pi_t = (A_1, X_1, A_2, X_2, \dots, A_{t-1}, X_{t-1}) \xrightarrow{\text{time}} P_t$$

action at time 1 reward of selected arm at time 1

$P_t \in \mathbb{P}_{k-1}$ probability simplex over
a probability dist. k dimensions
over $\{1, 2, \dots, k\}$. i.e., $q \in \mathbb{P}_{k-1}$

$$X_s = \sum_{j=1}^k x_{sj} \chi_{\{A_s=j\}}$$

$$\Rightarrow q = (q^{(1)}, \dots, q^{(k)})$$

with $q^{(j)} \geq 0$
 $\sum_{j=1}^k q^{(j)} = 1$.

- P_t is a function of all observations of actions and reward / loss up to and including time $(t-1)$.

- The adversary knows the selection function, i.e., adversary knows P_t before choosing (x_{t1}, \dots, x_{tk}) at time t .

(b) Player has access to secret random bits not known to adversary. Player uses these random bits to choose action according to policy $P_t = (P_{t1} P_{t2} \dots P_{tK})$, i.e.,

$$P(A_t=j | A_1, x_1, \dots, A_{t-1}, x_{t-1}) = P_{tj}$$

Thus adversary knows player policy but does not know player action, when deciding arm rewards.

③ Bandit feedback: Player gets to observe the reward / loss according to the specific action at time t , i.e.,

if $A_t=j$, then player gets to only observe x_{tj} (or y_{tj}).

(Other arms' rewards are NOT revealed).

Regret: In the adversarial setting, regret is with respect to the best fixed arm policy in hindsight.

$$R_n(\pi, x) = \left(\max_{1 \leq j \leq k} \sum_{t=1}^n x_{tj} \right) - E \left[\sum_{t=1}^n x_{tA_t} \right]$$

policy chosen by player
 reward
 environment x

= sequence $\{x_{tj}, j=1, \dots, k\}$
 $t=1, \dots, n$
 chosen by adversary

highest fixed arm reward in hindsight
 with $E[\cdot]$ over randomization
 by player.

Equivalently,

$$R_n(\pi, y) = E \left[\sum_{t=1}^n y_{tA_t} \right] - \left(\min_{1 \leq j \leq k} \sum_{t=1}^n y_{tj} \right)$$

loss environment y
 player's loss
 fixed arm with smallest loss in hindsight.

Aside: Why compare to best fixed arm in hindsight? One motivation is in the setting of classification, where the goal is to use n labeled samples, and return a weight vector characterizing a classifier.

Consider an algorithm that optimizes for the ERM (Empirical Risk Minimization) criterion.

It can be shown that developing an online learning algorithm that optimizes with the "best fixed arm in hindsight" criterion automatically leads to tight generalization bounds.

We will discuss this formally later when studying OCO (Online Convex Optimization).

Ref: Chap 5 — Online-to-Batch Conversions from the monograph: Online Learning and Online Convex Optimization, by Shai Shalev-Shwartz, NOW Publishers, Vol. 4, No. 2, 2011.

Exp 3 Algorithm :



Exponential weight algorithm for Explore and Exploit.

Idea 1: Suppose player had access to the rewards for all arms. Then, exponentially boost the probability of that arm

that has been best in hindsight at time t .

i.e., suppose all $\{y_{sj}, s=1, 2, \dots, t, j=1, \dots, k\}$

were revealed. Then, cumulative loss from arm j is $L_{t-1,j} = \sum_{s=1}^{t-1} y_{sj}$

Thus, choose arm j at time t

$$\alpha e^{-\eta L_{t-1,j}}$$

(or equivalently, $P_{tj} \propto e^{+\eta S_{t-1,j}}$)

cumulative reward until time ($t-1$)

Idea 2: We do not know $\{y_{tj}, j=1, \dots, k\}$ because of bandit feedback.

Thus, use Idea 1 (exp boost best in hindsight) using conditionally unbiased estimators the unknown y_{tj} .

(Seems strange, as we have no model for y_{tj} , so how can we create an estimator?)

Importance Sampling Estimators for $\{x_{tj}, y_{tj}\}$.

Let

$$\text{reward of arm } A_1 \quad x_s = \sum_{j=1}^k x_{sj} \chi_{\{A_s=j\}}$$

$$I_{t-1} = (A_1, x_1, A_2, x_2, \dots, A_{t-1}, x_{t-1})$$

or

$$= (A_1, \gamma_1, A_2, \gamma_2, \dots, A_{t-1}, \gamma_{t-1})$$

$$\gamma_s = \sum_{j=1}^k y_{sj} \chi_{\{A_s=j\}}.$$

(we are misusing the same notation for both as

$y_t = (1 - x_t)$, so conditioning on either is equivalent, and is usually clear from context which we are using).

$$P_{tj} = P(A_t=j \mid I_{t-1}) \quad \text{r.o.}$$

(the algo. we consider satisfies $P_{tj} > 0$ a.s.,
 $\forall j = 1, 2, \dots, k, t = 1, 2, \dots, n$)

The estimators:-

$$\hat{x}_{tj} = \frac{\chi_{\{A_t=j\}} \cdot x_t}{P_{tj}}$$

$$\hat{y}_{tj} = \frac{\chi_{\{A_t=j\}} \gamma_t}{P_{tj}}$$

Now, $\rightarrow E[\cdot]$ wrt to the randomness of the player's actions

$$E\left[\hat{x}_{tj} \mid I_{t-1}\right] = E\left[\frac{X_{\{A_t=j\}} x_t}{P_{tj}} \mid I_{t-1}\right]$$

because
P_{tj} is
a known
fn of I_{t-1}

$$= \frac{1}{P_{tj}} E\left[X_{\{A_t=j\}} \hat{x}_t \mid I_{t-1}\right] \leftarrow$$

Next, observe that given I_{t-1}, the action A_t and reward vector x_t are independent. This is because the adversary has access only to the player's policy (a deterministic function of I_{t-1}), but not the random bits the player uses to choose the action.

$$\therefore E\left[\hat{x}_{tj} \mid I_{t-1}\right] = \frac{1}{P_{tj}} E\left[X_{\{A_t=j\}} \mid I_{t-1}\right] x_{tj}$$

$\stackrel{= P_{tj}}{\sim}$
 ↴
 rand. reward
 chosen by adv.

$$= x_{tj}$$

Similarly, $E\left[\hat{y}_{tj} \mid I_{t-1}\right] = y_{tj}$

Aside: Consider the conditional variances of these estimators.

$$V_t [\hat{x}_{tj}^2] = E \left[\hat{x}_{tj}^2 \mid I_{t-1} \right] - x_{tj}^2$$

$$= E \left[\frac{x_{\{A_t=j\}} x_t^2}{P_{tj}^2} \mid I_{t-1} \right] - x_{tj}^2$$

$$= \frac{x_{tj}^2 (1 - P_{tj})}{P_{tj}} \quad \begin{matrix} \curvearrowright \text{good if} \\ \text{rewards are small} \end{matrix}$$

Similarly,

$$V_t [\hat{y}_{tj}^2] = \frac{y_{tj}^2 (1 - P_{tj})}{P_{tj}} \quad \begin{matrix} \curvearrowright \text{good if} \\ \text{losses are small} \end{matrix}$$

\downarrow
we will
use this
in the proofs

The Exp 3 Algorithm

(Sec. 11.3 in
text)

The discussion below is based on losses. Please read text for the one based on rewards.

Let $\hat{L}_{t,j} = \sum_{s=1}^t \hat{\gamma}_{sj}$, $j=1, 2, \dots, k$
 $t=1, 2, 3, \dots, n$
and $\hat{L}_{0,j} = 0$

$$\hat{\gamma}_{sj} = \frac{\chi_{\{A_s=j\}} \gamma_s}{P_{sj}} = \frac{\chi_{\{A_s=j\}} (1-\gamma_s)}{P_{sj}}$$

Algorithm input: $\eta > 0$

for each time $t=1, 2, \dots, n$:

- $P_{tj} = \left(\frac{e^{-\eta \hat{L}_{t-1,j}}}{\sum_{r=1}^k e^{-\eta \hat{L}_{t-1,r}}} \right)$ $\rightarrow j=1, 2, \dots, k$
 \downarrow
arms.

- Sample $A_t \sim P_{tj}$ $\left(\begin{array}{l} \text{using indep / secret} \\ \text{randomness} \end{array} \right)$

- Update $\hat{L}_{t,j} = \hat{L}_{t-1,j} + \hat{\gamma}_{tj}$

Thm (11.1 in text for rewards; below mirrors the proof for loss).

Setting as above. Then for $\eta = \sqrt{\frac{\ln k}{nk}}$,

$$R_n(r, y) \leq 2 \sqrt{nk \ln(k)}$$

Prof: Fix any arm i , and we will analyze the regret with respect to this fixed arm. As i is arbitrary, the analysis bounds the regret wrt any fixed arm (note this is diff. from regret decomposition we have been studying for stochastic bandits).

$$\rightarrow R_{ni} = E \left[\sum_{t=1}^n \gamma_t \right] - \underbrace{\sum_{t=1}^n y_{ti}}$$

Proof has 3 main steps, marked A, B, C below.

----- A: Express Regret wrt Estimates -----

$$\begin{aligned} E[\hat{I}_{ni}] &= E \left[\sum_{s=1}^n \hat{\gamma}_{si} \right] \\ &= E \left[\sum_{s=1}^n E[\hat{\gamma}_{si}] I_{s-1} \right] \end{aligned}$$

$$E[\hat{Z}_{ni}] = E\left[\sum_{s=1}^n y_{si}\right] = \sum_{s=1}^n y_{si} \rightarrow \textcircled{1}$$

$$E\left[y_t | I_{t-1}\right] = \sum_{j=1}^k P_{tj} y_{tj}$$

$$= \sum_{j=1}^k P_{tj} E\left[\hat{Y}_{tj} | I_{t-1}\right]$$

$\rightarrow \textcircled{2}$

Next, let

$$\hat{Z}_n = \sum_{t=1}^n \sum_{j=1}^k P_{tj} \hat{Y}_{tj}$$

$$\text{Then } E[\hat{Z}_n] = E\left[\sum_{t=1}^n E\left[\sum_{j=1}^k P_{tj} \hat{Y}_{tj} | I_{t-1}\right]\right]$$

$$\stackrel{\textcircled{2}}{\Rightarrow} = E\left[\sum_{t=1}^n E\left[y_t | I_{t-1}\right]\right]$$

(Tower rule)

$$\text{used in several steps} \rightarrow = E\left[\sum_{t=1}^n y_t\right] \rightarrow \textcircled{3}$$

$$\text{Thus, } R_{ni} = E\left[\sum_{t=1}^n y_t\right] - \sum_{t=1}^n y_{ti}$$

$$\therefore R_{ni} \stackrel{(3)(1)}{=} E[\hat{L}_n] - E[\hat{L}_{ni}] \rightarrow A \rightarrow ④$$

(Note: We will use $\hat{L}_{n,i} = \hat{L}_{ni}$ interchangeably)

----- Step (B): Track $(\hat{L}_n - \hat{L}_{ni})$ -----

$$\text{Let } W_t = e^{\eta t} \cdot \sum_{j=1}^k e^{-\eta \hat{L}_{tj}}, \quad t = 0, 1, 2, \dots, n$$

$$\text{with } W_0 = k \quad (\because \hat{L}_{0j} = 0 \text{ by definition}).$$

$$e^{\eta n} e^{-\eta \hat{L}_{ni}} \leq \sum_{j=1}^k e^{-\eta \hat{L}_{nj}} e^{\eta n}$$

$$= W_n = W_0 \left(\frac{W_1}{W_0} \right) \cdots \left(\frac{W_n}{W_{n-1}} \right)$$

→ ⑤

$$= k \prod_{t=1}^n \left(\frac{W_t}{W_{t-1}} \right)$$

$$\text{Now, } \left(\frac{W_t}{W_{t-1}} \right) = \frac{e^{\eta t} \sum_{j=1}^k e^{-\eta \hat{L}_{tj}}}{e^{\eta(t-1)} \sum_{j=1}^k e^{-\eta \hat{L}_{t-1,j}}} = e^{\eta} \sum_{j=1}^k \frac{e^{-\eta \hat{L}_{t-1,j}} \cdot e^{-\eta \hat{L}_{tj}}}{\sum_{r=1}^k e^{-\eta \hat{L}_{t-1,r}}}$$

$$= \sum_{j=1}^k P_{tj} e^{\eta (1 - \hat{L}_{tj})}$$

Now, observe that:

$$\textcircled{a} \quad \text{For } x \leq 1, (1+x) \stackrel{(i)}{\leq} e^x \stackrel{(ii)}{\leq} 1+x+x^2.$$

$$\textcircled{b} \quad \hat{\gamma}_{tj} = \frac{x_{\{A_t=j\}} \gamma_t}{P_{tj}} \stackrel{\text{e [0,1]}}{=} \begin{cases} 0 & \text{if } A_t \neq j \\ \frac{\gamma_t}{P_{tj}} & \text{if } A_t = j \end{cases}$$

$0 < \hat{\gamma}_{tj} < 1$

$$\Rightarrow (1 - \hat{\gamma}_{tj}) \leq 1$$

$$\therefore \frac{w_t}{w_{t-1}} = \sum_{j=1}^k P_{tj} \cdot e^{\eta(1 - \hat{\gamma}_{tj})}$$

$$\stackrel{(ii)}{\leq} \sum_{j=1}^k P_{tj} \cdot \left(1 + \eta(1 - \hat{\gamma}_{tj}) + \eta^2 (1 - \hat{\gamma}_{tj})^2 \right)$$

$$= \left(1 + \eta \sum_{j=1}^k P_{tj} (1 - \hat{\gamma}_{tj}) + \eta^2 \sum_{j=1}^k P_{tj} (1 - \hat{\gamma}_{tj})^2 \right)$$

$$\stackrel{(i)}{\leq} e^{\eta \sum_{j=1}^k P_{tj} (1 - \hat{\gamma}_{tj}) + \eta^2 \sum_{j=1}^k P_{tj} (1 - \hat{\gamma}_{tj})^2}$$

$$\therefore e^{n\eta} \cdot e^{-\eta \hat{L}_{ni}} \leq \prod_{t=1}^n \left(\frac{w_t}{w_{t-1}} \right) \quad (\text{from } ⑤)$$

$$\therefore \ln k + \eta \sum_{t=1}^n \sum_{j=1}^k p_{tj} (1 - \hat{\gamma}_{tj}) + \eta^2 \sum_{t=1}^n \sum_{j=1}^k p_{tj} (1 - \hat{\gamma}_{tj})^2 \leq \ell$$

$$\therefore n - \hat{L}_n \quad (\text{see below})$$

$$\therefore (n - \hat{L}_{ni}) \leq \frac{\ln k}{\eta} + \sum_{t=1}^n \sum_{j=1}^k p_{tj} (1 - \hat{\gamma}_{tj})$$

$\log(\cdot)$ and
divide by η .

$$+ \eta \sum_{t=1}^n \sum_{j=1}^k p_{tj} (1 - \hat{\gamma}_{tj})^2$$

$$\text{Recall: } \hat{L}_n = \sum_{t=1}^n \sum_{j=1}^k p_{tj} \hat{\gamma}_{tj}, \text{ and note}$$

$$\text{that } \sum_{t=1}^n \sum_{j=1}^k p_{tj} = n$$

→ ⑥

$$\Rightarrow (\hat{L}_n - \hat{L}_{ni}) \leq \frac{\ln k}{\eta} + \eta \sum_{t=1}^n \sum_{j=1}^k p_{tj} (1 - \hat{\gamma}_{tj})^2$$

$E[\cdot] = R_{ni}$

→ ⑥

(C): Bound the expectation of RHS in ⑥:

$$E \left[\sum_{t=1}^n \sum_{j=1}^k p_{tj} (1 - \hat{y}_{tj})^2 \right]$$

$$= E \left[\sum_{t=1}^n \sum_{j=1}^k p_{tj} \left(1 - \frac{\chi_{\{A_t=j\}} y_t}{p_{tj}} \right)^2 \right]$$

$$= \sum_{t=1}^n E \left[\sum_{j=1}^k p_{tj} \left(1 + \frac{\chi_{\{A_t=j\}} y_t^2}{p_{tj}^2} - \frac{2 \chi_{\{A_t=j\}} y_t}{p_{tj}} \right) \right]$$

$$= \sum_{t=1}^n E \left[1 + \sum_{j=1}^k E \left[\frac{\chi_{\{A_t=j\}} y_t^2}{p_{tj}} - 2 \chi_{\{A_t=j\}} y_t \mid I_{t-1} \right] \right]$$

Note: $\chi_{\{A_t=j\}} y_t = \chi_{\{A_t=j\}} y_{tj}$

$$= \sum_{t=1}^n E \left[1 + E \left[\sum_{j=1}^k \frac{\chi_{\{A_t=j\}} y_{tj}^2}{p_{tj}} \mid I_{t-1} \right] \right]$$

$$- 2 E \left[\sum_{j=1}^k \chi_{\{A_t=j\}} y_{tj} \mid I_{t-1} \right]$$

$\underbrace{\quad}_{= y_t}$

$\therefore E[\chi_{\{A_t=j\}} \mid I_{t-1}] = p_{tj}$

$$= \sum_{t=1}^n E \left[1 - 2\gamma_t + E \left[\sum_{j=1}^k \frac{\chi_{\{A_t=j\}} y_{tj}^2}{P_{tj}} \mid \mathcal{I}_{t-1} \right] \right]$$

$$= \sum_{t=1}^n E \left[1 - 2\gamma_t + \sum_{j=1}^k y_{tj}^2 \right]$$

$$= \sum_{t=1}^n E \left[1 - 2\gamma_t + \gamma_b^2 - \gamma_t^2 + \sum_{j=1}^k y_{tj}^2 \right]$$

$$= \sum_{t=1}^n E \left[\underbrace{(1-\gamma_t)^2}_{\leq 1} + \underbrace{\sum_{\substack{j=1 \\ j \neq A_t}}^k y_{tj}^2}_{\leq (k-1)} \right] \leq nk$$

\therefore

$$\boxed{E \left[\sum_{t=1}^n \sum_{j=1}^k P_{tj} (1 - \hat{\gamma}_{tj})^2 \right] \leq nk} \rightarrow \textcircled{C}$$

Substituting in $\textcircled{6}$ and taking expectations,

$$R_{ni} \leq \frac{\ln k}{\eta} + \eta nk = 2 \sqrt{n k \ln(k)}$$

$$\eta = \sqrt{\frac{\ln k}{nk}}$$

Comment: We can tighten the constant by $\sqrt{2}$,

$$\text{i.e., } R_n \leq \sqrt{2nk \ln(k)}$$

Please read text, Thm 11.2 for the proof.

Note: Instead of bandit feedback, suppose we had full feedback, i.e., the player gets to see rewards/losses for all arms (but incurs the one that the player chose to play).

A decision-theoretic
→ Hedge algorithms → generalization of online learning and an application to boosting, Freund and Schapire, 1997

$$\text{Then regret} \leq \sqrt{\frac{n \ln(k)}{2}}$$

* No \sqrt{k} term! → The gain due to full information

(See Note 3, Chap 11 in text for discussion.)

To see this scaling, recall (5) :

$$\underbrace{(\hat{L}_n - \hat{L}_{ni})}_{E[\cdot] = R_{ni}} \leq \frac{\ln k}{\eta} + \eta \sum_{t=1}^n \sum_{j=1}^k P_{tj} (1 - \hat{y}_{tj})^2$$

Suppose $\hat{y}_{tj} = y_{tj}$ (i.e., no noise in the estimator because of full information)

Then

$$\sum_{t=1}^n \sum_{j=1}^k P_{tj} (1 - \hat{y}_{tj})^2 = \sum_{t=1}^n \sum_{j=1}^k P_{tj} (1 - y_{tj})^2 \underbrace{\leq 1}_{\leq 1} \leq n$$

(i.e., bound is $\left(\frac{\ln k}{\eta} + \eta n \right)$. Choose

$$\eta = \sqrt{\frac{\ln k}{n}} \text{ to balance terms} \Rightarrow$$

$$R_n \leq 2 \sqrt{n \ln k} \quad \left(\begin{array}{l} \text{Analysis can be} \\ \text{sharpened to} \\ \text{provide tighter constant} \end{array} \right)$$