



# Adversarial bandits:

15 | 4 | 20

## k-armed adv. bandit:

$(x_t)_{t=1}^n$  is a seq of reward vectors  
with  $x_t \in [0, 1]^k$ ;

In each round the learner chooses a  
distribution over actions  $P_t \in \mathcal{P}_{k-1}$ .  
Then action  $A_t \in [k]$  is sampled from  
 $P_t$ , and learner gets  $x_{tA_t}$ ;

## Alg:

Adv selects  $(x_t)_{t=1}^n$   $x_t \in [0, 1]^k$

for rounds  $t = 1, 2, \dots, n$

learner selects  $P_t \in \mathcal{P}_{k-1}$ , and samples  
 $A_t$  from  $P_t$ ;

learner sees  $x_{tA_t}$

Policy:  $\pi: ([k] \times [0,1])^* \rightarrow \mathcal{P}_{k-1}$  }

mapping history sequences to distributions  
over actions;

$A_1, X_1, A_2, X_2, A_3, \dots$

$$R_n(\pi, \alpha) = \max_{i \in [k]} \sum_{t=1}^n r_{ti} - \mathbb{E} \left[ \sum_{t=1}^n r_{tA_t} \right]$$

Comparing against a fixed action in  
highlight;

= Randomness - learner;

- Action chosen in round  $t$  may depend  
upon actions chosen at  $s < t$ , and  
the observed rewards.

$$R_n^* (\pi) = \sup_{\pi \in \Sigma^{n \times h}} R_n (\pi, \pi)$$

If learner is deterministic:

linear regret can be forced;

look at the algorithm the learner is using;

If  $I_t = 1$ , let the rewards be

0	1	1	-	-	...
---	---	---	---	---	-----

If  $I_t \neq 1$  let rewards be

1	0	0	-	-	...
---	---	---	---	---	-----

• Now if  $\{t \mid I_t = 1\} \geq n/2$ , consider playing arm 2 all the time;

• Report<sup>2</sup>  $\geq \underline{1/2}$ .

• of  $\{t \mid I_t = 1\} \leq n/2$ ;  
consider playing 1 always;  
∴  $\{t \mid I_t \neq 1\} \geq \underline{1/2}$  ∴ get a report of  $\underline{1/2}$   
here,

• Focused to look at randomized learners: }

Regret is now a random variable,

- we need bounds in high probability or in expectation on  $R_n$ .

- Strategies we've seen is for one det.

---

• Randomization is in the rewards!

So none will work well in the adv. setting

- Converse<sup>2</sup>

Will <sup>an</sup> adversarial bandit strategy work in the stochastic setting?

- Small expected regret in stochastic setting.<sup>2</sup>

• Let  $\pi$  be an adv bandit strategy, and  $\nu = (\nu_1, \dots, \nu_k)$  a stochastic bandit,

$\text{Supp}(\nu_i) \subseteq [0, 1] \forall i$ ;

- Let  $X_{t,i}$  be sampled from  $\nu_i$  for  $i \in [k]$   
and  $t \in [n]$ ;  
Assume  $X_{t,i}$  are mutually indep

$$R_n(\pi, \mathcal{X}) = \max_{i \in [k]} \mathbb{E} \left[ \sum x_{ti} - x_{tAT} \right]$$

$$\leq \mathbb{E} \left[ \max_{i \in [k]} \sum x_{ti} - x_{tAT} \right]$$

↑ Stoch syst.

$$= \mathbb{E} \left[ R_n(\pi, \mathcal{X}) \right] \leq R_n^*(\pi).$$

↑ Adv regret!

- Worst case Stochastic regret is upper bounded by worst case adversarial regret.

- Key idea: A mechanism for estimating rewards of unplayed arms.

Now  $P_t$  is the conditional distribution of actions played in round  $t$

$$\therefore P_{ti} = P[A_t = i | A_{(1)}, X_{(1)}, \dots, X_{t-1}]$$



conditioned on history

Define:

importance-weighted estimator of  $\pi_{ti}$ :

$$\hat{X}_{ti} = \frac{\mathbb{1}\{A_t = i\} X_t}{P_{ti}}$$

Let  $E_t[\cdot] = E[\cdot | A_{(1)}, \dots, X_{t-1}]$  conditional exp given history.

If  $A_{ti} = \mathbb{1}\{A_t = i\}$  then clearly.

$$X_t A_{ti} = \pi_{ti} A_{ti}$$

$$\therefore \hat{X}_{ti} = \frac{A_{ti} X_t}{P_{ti}} = \frac{\pi_{ti} A_{ti}}{P_{ti}}$$



$$\therefore \mathbb{E}_t[\hat{X}_{ti}] = \mathbb{E}_t \left[ \frac{\alpha_{ti} A_{ti}}{P_{ti}} \right]$$

$$= \mathbb{E}_t \left[ \frac{\alpha_{ti} A_{ti}}{P_{ti}} \mid A_{1i}, X_{1i}, \dots, X_{t-1i} \right]$$

$$= \frac{\alpha_{ti}}{P_{ti}} \mathbb{E}_t [A_{ti} \mid \dots]$$

$$= \frac{\alpha_{ti}}{P_{ti}} \cdot [1 \cdot P_{ti} + 0 \cdot (1 - P_{ti})]$$

$$= \alpha_{ti}$$

$\therefore \hat{X}_{ti}$  - an unbiased est. of  $\alpha_{ti}$ .

• What about variance?

$$\text{Now: } \underline{V_t[u]} = \mathbb{E}_t \left[ (u - \mathbb{E}_t(u))^2 \right]$$

$$\begin{aligned}
 \therefore V_t[\hat{X}_{ti}] &= E_t[\hat{X}_{ti}^2] - (E_t[\hat{X}_{ti}])^2 \\
 &= E_t\left[\frac{A_{ti} x_{ti}^2}{P_{ti}}\right] - x_{ti}^2 \\
 &= \frac{x_{ti}^2}{P_{ti}} - x_{ti}^2 = \frac{x_{ti}^2 (1 - P_{ti})}{P_{ti}}
 \end{aligned}$$

= this can be large if  $P_{ti}$  is small  
and  $x_{ti} > \epsilon$

Another estimator:

$$\hat{X}_{ti} = 1 - \frac{\mathbb{1}\{A_t = i\} (1 - X_t)}{P_{ti}}$$

This too is unbiased;

Using  $y_{ti} = 1 - x_{ti}$ ,  $Y_t = 1 - X_t$

&  $\hat{Y}_{ti} = 1 - \hat{X}_{ti}$  we get

$$Y_{ti}^A = \frac{\mathbb{1}\{A_t=i\} Y_t}{P_{ti}} \quad \uparrow$$

← again:  $E[Y_{ti}^A] = 1 - p_{ti}$

Loss based estimator!

$$V_t[Y_{ti}^A] = y_{ti}^2 \left( \frac{1 - P_{ti}}{P_{ti}} \right)$$

∴ which is better?   
 smaller reward,  $X_{ti}^A$  for arm  $i$   $\Rightarrow$  better

Exp3 algorithm: Exp wts. alg for Exp & Exploitation.

Let  $S_{ti}^A = \sum_{s=1}^t X_{si}^A$ ,  $X_{ti}^A = 1 - Y_{ti}^A$

Clearly  $V_t[X_{ti}^A] = V_t[Y_{ti}^A]$

Normal to play action with large estimated rewards with by softmax probability.

Can map  $\hat{Q}_{t,i}$  — and probability;

$$P_{t,i} = \frac{\exp(\eta \hat{Q}_{t-1,i})}{\sum_{j=1}^k \exp(\eta \hat{Q}_{t-1,j})}$$

$\eta$  - learning rate;

Use this to get EXP3 algorithm.

1) Input:  $n, k, \eta$

2)  $\hat{s}_{0i} = 0 \forall i$

3) for  $t = 1, \dots, n$  do

4) 
$$P_{ti} = \frac{\exp(\eta \hat{s}_{t-1, i})}{\sum_{j=1}^k \exp(\eta \hat{s}_{t-1, j})}$$

5) Sample  $A_t \sim P_t$  & observe  $X_t$

$$\textcircled{6} \quad \hat{s}_{ti} = \hat{s}_{t-1, i} + \frac{\mathbb{1}\{A_t = i\} (1 - X_t)}{P_{ti}}$$

⑦ end for;

Then:  $\alpha \in [0, 1]^{n \times k}$ , set  $\eta = \sqrt{\log(k) / (nk)}$ .

Then  $R_n(\Pi, \alpha) \leq 2 \sqrt{nk \log k}$ .