# Week 3

## Lecture 13

The expected regret can also be:

$$E\left[\sum_{t=1}^{n} x_{tI_t}\right] - E\left[\sum_{t=1}^{n} \min_i x_{ti}\right] \qquad (2a)$$

$\underbrace{\phantom{E\left[\sum_{t=1}^{n} \min_i x_{ti}\right]}}$

Here the benchmark is that we choose the best action in each round

$$E\left[\sum_{t=1}^{n} x_{tI_t}\right] - \min_i E\left[\sum_{t=1}^{n} x_{ti}\right] \qquad (2)$$

$\underbrace{\phantom{\min_i E\left[\sum_{t=1}^{n} x_{ti}\right]}}$

Here the benchmark is the <u>one</u> arm which if played for all rounds will give min loss

Though eq 2a is a better criterion, but it is harder to achieve and hence expected regret is defined by eq 2

As long as $x_t$ are generated randomly the regret generated by eq 2a $\geq$ regret by eq 2

This is because:

$$E\left[\sum_{t=1}^{n} \min_i x_{ti}\right] \leq \min_i E\left[\sum_{t=1}^{n} x_{ti}\right]$$

If $x_t$ are generated according to some sequence then eq 2 becomes:

$$R(n, \pi) = \mathbb{E}\left[\sum_{t=1}^{n} x_{tI_t}\right] - \min_{i}\left[\sum_{t=1}^{n} x_{ti}\right] \qquad \text{—②b}$$

In eq 2, the first term has 2 levels of randomness : → randomness of loss by the adversary

→ ——do—— arm selection by learner

$$\mathbb{E}\left[R(n, \pi)\right] = \mathbb{E}\left[\sum_{t=1}^{n} x_{tI_t}\right] - \min_{i} \mathbb{E}\left[\sum_{t=1}^{n} x_{ti}\right] \qquad \text{—②}$$

↑

Psuedo regret, also known as

The learner chooses the next action based on history, hence the expected regret should ideally be over a particular sequence.

$$\therefore \quad \mathbb{E}\left[R(n, \pi, \{x_t\})\right] = \mathbb{E}\left[\sum_{t=1}^{n} x_{tI_t}\right] - \min_{i}\left[\sum_{t=1}^{n} x_{ti}\right] \qquad \text{—②c}$$

Assuming $\{X_t\}$ to be stochastic, $\mathbb{E}$ have to be an both the terms an the R.H.S. of eq 2...

Hence revisiting the definitions again :

Expected Regret :
$$\mathbb{E}\left[R(n, \pi)\right] = \mathbb{E}\left[\sum_{t=1}^{n} X_{tI_t}\right] - \mathbb{E}\left[\sum_{t=1}^{n} \min_{i} X_{ti}\right]$$
$$\geq \mathbb{E}\left[\sum_{t=1}^{n} X_{tI_t}\right] - \min_{i} \mathbb{E}\left[\sum_{t=1}^{n} X_{ti}\right]$$

↑

# Psuedo regret

Whatever has been discussed till now about bandits is the <u>Adversarial Bandit setting</u>

[ starting from Lec 12 ]

## Importance Sampling

At round $t$

$$P_t = (P_{t1}, P_{t2}, \dots P_{tk})$$

Learner chooses $I_t \sim P_t$

Let $I_t = i$,

Learner suffers loss $x_{ti}$

At round $t$ we can define an estimator for all $\{x_t\}$

Importance weighted estimator :

$$\forall j \in [K] \quad \hat{x}_{tj} = \frac{x_{tj}}{P_{tj}} 1\{I_t=j\} = \begin{cases} \frac{x_{ti}}{P_{ti}} & \text{for } j=i \\ \\ 0 & \text{for } j \neq i \end{cases}$$

$$\mathbb{E}[\hat{x}_{tj}] = \sum_{i=1}^{K} \frac{x_{tj}}{P_{tj}} 1\{i=j\} \times P_{ti}$$

$$= \frac{x_{tj}}{P_{tj}} \times P_{tj}$$

$$= x_{tj}$$

Hence the importance weighted estimator is an unbiased estimator of the losses.

Hence though we don't have information for all the arms/actions, but we have unbiased estimators for all losses, which we can use to update the weights.

---