

Bandits with full feedback

↳ partial feedback (only get reward for chosen action)

Full feedback - observe rewards / costs of all actions (arms), not only chosen one

e.g. expert feedback (Online learning with experts)
stock market

Problem setup (full feedback)

K actions / arms, T rounds

At each round t

adversary chooses $c_t(a)$ $\forall a = 1 \dots K$
algorithm picks a_t and incur $c_t(a_t)$
costs of all $\{c_t(a)\}$ revealed

Goal: Do as well as best action.

Online learning with experts

K experts, T rounds

adversary chooses observation x_t and label z_t
Reveals x_t only, expert prediction $\hat{z}_{i,t} \quad i = 1 \dots K$
Algorithm picks expert $e \in \{1 \dots K\}$
Incur loss $c_t(e) = c(\hat{z}_{e,t}, z_t)$ on seeing true label z_t .
known
Costs for all experts revealed

Note: Known costs for bandits for easy.

Note: Stochastic costs for full feedback is easy.

Play lowest average cost action

$$R(T) \leq \sqrt{\log K T} \left(1 + \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{3}} + \dots \right)$$

$$O\left(\sqrt{\frac{\log K T}{n_t(a)}}\right)$$

$$\frac{1}{\sqrt{t}} \leq \frac{2}{\sqrt{t} + \sqrt{t-1}} = 2(\sqrt{t} - \sqrt{t-1})$$

$$\sum \frac{1}{\sqrt{t}} = 2\sqrt{T}$$

$$R(T) = O(\sqrt{T \log K T})$$

vs. Bandit $O(\sqrt{KT \log KT})$

Adversarial setting

Oblivious or Adaptive adversary
 costs don't) 'do' depend on algorithm's choices.

(Cumulative) Regret $R(T) = \sum_{t=1}^T c_t(a_t) - \sum_{t=1}^T c_t(a^*)$

 $a^* = \arg \min_a \sum_{t=1}^T c_t(a)$ - best in hindsight action.

Pseudo-Regret $\sum_{t=1}^T c_t(a_t) - \min_a E\left[\sum_{t=1}^T c_t(a)\right]$ reduces to iid best in foresight action know dist'g of costs

For deterministic adversary → Regret
 " randomized " → Regret + Pseudo-regret (weaker)

Note: Pseudo-regret cannot exceed regret.

Note: For iid costs/rewards, \sqrt{T} pseudo-regret extend to regrett $\log T$ do not . . .

Binary prediction $\hat{z}_{i,t} \in \{\pm 1\}$

Majority vote of experts who have not made mistake in past.

Thm: If a perfect expert exists, then majority vote makes at most $\log K$ mistakes.

Prog: S_t - experts who make no mistake upto t.

$$w_t = |S_t| , w_1 = K \leftarrow \dots \subset S_1 \subset \dots \subset S_t$$

$w_t \geq 1$ & t bar of perfect expert $\leftarrow t \sim \dots$
 If algo makes a mistake at t, $w_{t+1} \leq w_t/2 \leftarrow$
 lower majority of experts
 in S_t are wrong.

Sketch idea:

w_t - notion of \hat{w}_t weight on experts
 $w_1 \leq \square$ w_t doesn't increase over time
 $w_T \geq \square \leftarrow$ some notion of $\text{cost}(a^*)$

What if \hat{w}_t perfect expert?

Weighted Majority Vote

$$w_t(a) = 1 + a$$

For each t
 Make prediction using weighted \hat{w}_t maj vote.

For each expert i

$$w_i \leftarrow w_i \text{ if correct}$$

$$w_i \leftarrow (1-\epsilon) w_i \text{ if expert incorrect.}$$

Thm: # mistake by weighted maj vote $\leq \frac{2}{1-\epsilon} \text{cost}^* + \frac{2}{\epsilon} \ln K$

$$\text{cost}^* = \underbrace{\sum_{t=1}^T \text{cost}(a^*)}_{0} \downarrow \ln K$$

Beyond binary prediction

Note: Any deterministic algo. has total cost T even for oblivious adversary.

Hedge algorithm $\epsilon \in (0, \frac{1}{2})$

$$w_t(a) = 1 + a$$

$$\alpha_t = \frac{w_t(a)}{\sum_{a' \in A} w_t(a')}$$

→ Sample arm/expert at from distribution $\pi_t = \sum_a w_t(a)$
 observe costs of all arms
 $w_{t+1}(a) = w_t(a) (1-\varepsilon)^{c_t(a)}$ \Leftrightarrow $\begin{cases} 1 & c_t(a) = 0 \text{ for binary prediction} \\ 1-\varepsilon & c_t(a) = 1 \end{cases}$

Thm: Bounded costs ≤ 1 . Consider adaptive adversary s.t. $\sum_t c_t(a^*) \leq uT$ for some known u , then Hedge with $\varepsilon = \sqrt{\frac{\ln K}{uT}}$ satisfies.

$$\overline{E\left[\sum_{t=1}^T c_t(a_t) - \sum_{t=1}^T c_t(a^*)\right]} \leq \frac{2\sqrt{uT \ln K}}{w_t(a^*)}.$$

Proof: ① $w_t = \sum_a w_t(a)$
 $w_1 = \frac{K}{1}$, $w_{T+1} \geq (\text{cost}(a^*))$ ($w_T \geq 1$ for majority assuming 1 best expert exists)
 $w_{T+1} > w_{T+1}(a^*) = (1-\varepsilon) \frac{\sum_t c_t(a^*)}{w_t(a^*)} = (1-\varepsilon) \frac{\text{cost}(a^*)}{w_t(a^*)}$

$$② \quad \frac{w_{t+1}}{w_t} = \frac{\sum_a (1-\varepsilon)^{c_t(a)} w_t(a)}{\sum_a w_t(a)} p_t(a) \quad \left(\frac{w_{t+1}}{w_t} \leq \frac{1}{2} \text{ maj vote for binary prediction} \right)$$

$$\leq \sum_a (1-\varepsilon^{c_t(a)}) p_t(a) = 1 - \varepsilon \sum_a (c_t(a)) p_t(a)$$

$$(1-\varepsilon)^c \leq 1 - \varepsilon c \quad \text{if } c \leq 1 \quad (\text{bounded costs})$$

$$1 - \varepsilon \leq (1 - \varepsilon c)^{1/c} = 1 - \varepsilon + \dots$$

$$\ln \frac{w_{t+1}}{w_t} \leq \ln \underbrace{(1 - \varepsilon \sum_a c_t(a) p_t(a))}_{< -\varepsilon \sum_a c_t(a) p_t(a)} \quad \because 1 - x \leq e^{-x}$$

$$\sum_t \ln \frac{w_{t+1}}{w_t} = \ln \prod_t \frac{w_{t+1}}{w_t} = \ln \frac{w_{T+1}}{w_1}$$

$$\sum_t \varepsilon E[c_t(a_t)] = -\sum_t \ln \frac{w_{t+1}}{w_t} = -\ln \frac{w_{T+1}}{w_1} \leq -\ln \frac{(1-\varepsilon)^{\sum_t c_t(a^*)}}{K}$$

$$E[\sum_t c_t(a_t)] \leq \ln \frac{K}{1-\varepsilon} + \sum_t c_t(a^*) \frac{1}{\varepsilon} \ln \frac{1}{1-\varepsilon}$$

$$- \epsilon_t \cdot \varepsilon \quad \sim \leq 1 + \varepsilon \quad \text{if } \varepsilon \leq \frac{1}{2}$$

$$E \left[\sum_t c_t(a_t) - \sum_t c_t(a^*) \right] \leq \underbrace{\frac{\ln K}{\varepsilon}}_{\sum_t c_t(a^*)} + \varepsilon \sum_t c_t(a^*)$$

$$\sum_t c_t(a^*) = \ln K$$

$$\varepsilon = \sqrt{\frac{\sum_t c_t(a^*)}{\ln K}} = O\left(\sqrt{\frac{\sum_t c_t(a^*) \ln K}{nT}}\right)$$

2

① extend to unbounded costs

② Comparator class C_T - action sequence (a^*, a^*, \dots, a^*)

$$\text{Regret} \quad \sum_{t=1}^T c_t(a_t) - \min_{\substack{y_1, \dots, y_T \\ \in C_T}} \sum_{t=1}^T c_t(a_1, \dots, a_{t-1}, y_t) \leq \eta \quad \begin{matrix} \text{Oblivious} \\ \text{adversary} \end{matrix}$$

$$\text{Policy regret} \quad \sum_{t=1}^T c_t(a_{1:t}) - \min_{\substack{y_1, \dots, y_T \in C_T}} \sum_{t=1}^T c_t(y_{1:t}) = S(T) \quad \begin{matrix} \text{adaptive} \\ \text{adversary} \end{matrix}$$

↓
times (tallying
bandit)