

Recap Adversarial full-feedback setting \equiv Online learning with experts

N action/experts, T rounds

At each round

adversary picks $c_t(a) \forall a$

algorithm picks action/expert a_t

costs of all actions/expert $c_t(a) \forall a \equiv$ full feedback.

Hedge

$w_1(e) = 1$ for each expert

Sample expert e_t from distribution $p_t(e) = \frac{w_t(e)}{\sum_e w_t(e)}$

observe costs for all experts

multiplicative $w_{t+1}(e) \leftarrow w_t(e) (1 - \epsilon)^{c_t(e)}$ $\epsilon \in (0, 1/2)$

$$\text{Regret}(\text{Hedge}) = O\left(\sqrt{\sum_{t=1}^T c_t(a_t^*)} \log N\right)$$

exponential $w_{t+1}(e) \leftarrow w_t(e) e^{-\epsilon c_t(e)}$

Adversarial Bandits - Multiarmed bandit with adversarially chosen costs for each arm.

Reduction to full feedback

create expert corresponding to each action.

\hookrightarrow always picks a specific action

Algorithm: At each round t

Use Hedge to draw an expert e_t from $p_t \propto w_t$

Use e_t to pick arm a_t (TBD) \leftarrow

Observe cost $c_t(a_t)$ - only for a_t , NOT $\forall a$.

Define fake costs $\hat{c}_t(e) \forall e$. (TBD) \leftarrow

Return to Hedge

More general problem:

Adversarial bandits with expert feedback

- each expert can pick a different arm in different rounds.

K arms/actions, N experts, T rounds

In each round $t \in [T]$

Adversary picks cost $c_t(a) \geq 0 \quad \forall a$

Each expert recommends arm $a_{t,e}$

Algorithm picks arm $a_t \in [K]$ & receives cost $c_t(a_t)$

Total cost of each expert $\text{cost}(e) = \sum_{t=1}^T c_t(a_{t,e})$

Goal: Minimize regret relative to best expert — (than best arm)

$$R(T) = \text{cost}(\text{Alg}) - \min_e \text{cost}(e) \\ = \sum_{t=1}^T c_t(a_t) - \min_e \sum_{t=1}^T c_t(a_{t,e})$$

Note: If each expert recommends fixed arm then $N=K$
& $R(T) = \sum_{t=1}^T c_t(a_t) - \min_a \sum_{t=1}^T c_t(a) \leftarrow$ best arm in hindsight.

Assumptions: experts can't learn over time \Rightarrow all $a_{t,e}$ are chosen in advance.
focus on oblivious adversary \Rightarrow all costs are selected in advance.
bounded costs $c_t(a) \leq 1 \quad \forall t, a$.

Thm: Alg achieves $E[R(T)] = O(\sqrt{KT \ln N})$

lower bound = $\Omega(\sqrt{KT})$

$\Omega(\min(T, \sqrt{KT \ln N} / \ln K))$

Proof: $\hat{R}_{\text{Hedge}}(T) = \text{cost}(\text{Hedge}) - \min_e \text{cost}(e)$

$R_{\text{Hedge}}(T) = \text{cost}(\text{Hedge}) - \min_e \text{cost}(e)$

$\text{cost} = \sum_{t=1}^T c_t(\cdot)$

If $E[\hat{c}_t(e) | \mathcal{P}_t] = c_t(e) \quad \forall e$ then $E[R_{\text{Hedge}}(T)] \leq E[\hat{R}_{\text{Hedge}}(T)]$

Proof: $E[\hat{c}_t(e) | \mathcal{P}_t] = \sum_e \Pr(e_t = e | \mathcal{P}_t) E[\hat{c}_t(e_t) | \mathcal{P}_t]$

$$\begin{aligned}
&= \sum_e p_t(e) c_t(e) \\
&= E[c_t(e) | P_t] \\
\Rightarrow E[\hat{\text{Cost}}(\text{Hedge})] &= E[\text{Cost}(\text{Hedge})] \\
E[\min_e \text{cost}(e)] &\leq \min_e E[\text{cost}(e)] = \min_e E[\text{cost}(e)] \\
\Rightarrow E[R_{\text{Hedge}}(T)] &\leq E[\hat{R}_{\text{Hedge}}(T)]
\end{aligned}$$

① * Define fake costs.

Inverse Propensity Score

$$\hat{c}_t(a) = \begin{cases} c_t(a) / q_t(a) & a_t = a \\ 0 & \text{o.w.} \end{cases}$$

$$q_t(a) = \Pr(a_t = a | P_t)$$

Fake cost of each expert $\hat{c}_t(e) = \hat{c}_t(a_{t,e})$

$$\begin{aligned}
E[\hat{c}_t(e) | P_t] &= E[\hat{c}_t(a_{t,e}) | P_t] = \Pr(a_t = \underline{a_{t,e}} | P_t) \frac{c_t(a_{t,e})}{q_t(a_{t,e})} \\
&\quad + \Pr(a_t \neq a_{t,e} | P_t) \cdot 0 \\
&= c_t(a_{t,e}) = c_t(e).
\end{aligned}$$

② Define arm selection.

Hedge guarantee requires bounded^{fake} costs.

$\Rightarrow q_t(a)$ needs to be sufficiently large.

But also want to follow e_t to ensure low costs.

w.p. $1-\gamma$ follow $e_t \Rightarrow a_t = a_{t,e_t}$

γ choose a_t randomly. $\Rightarrow q_t(a) \geq \frac{\gamma}{K} \neq 0$.

Algorithm (EXP4) - exploration, exploitation, experts, exponentiation.

Inputs: N experts, K actions, ϵ for Hedge, γ for exploration.

At each round t

Call Hedge to draw an expert e_t using $P_t \propto w_t$

if $1-r$ follow expert e_t , $a_t = a_{t,e_t}$ }
 else pick a_t unij @ random

Observe $c_t(a_t)$

Define fake costs \hat{c}_t

$$\hat{c}_t(e) = \begin{cases} \frac{c_t(a_t)}{P(a_t = a_{t,e} | P_t)} & a_t = a_{t,e} \\ 0 & \text{o.w.} \end{cases}$$

Return fake costs to Hedge.

Note: fake costs are NOT oblivious.

but depend on observed costs of choices made by Alg in past.

need adaptive adversary guarantees for Hedge

$$E[\text{cost}(\text{Exp}_4)] \leq E[\text{cost}(\text{Hedge})] + rT$$

$$\Rightarrow E[\text{Regret}_{\text{Exp}_4}(T)] \leq E[R_{\text{Hedge}}(T)] + rT$$

$$\leq E[\hat{R}_{\text{Hedge}}(T)] + rT$$

Hedge bound $\sum_{t=1}^T \hat{c}_t(e^*) \leq \underline{uT}$ $\leftarrow \hat{c}_t(e) \leq \frac{1}{P_t(a_{t,e})} \leq \frac{K}{r}$

then $E[\hat{R}_{\text{Hedge}}(T)] \leq \sqrt{\underline{uT} \log N}$

$$\Rightarrow E[\text{Regret}_{\text{Exp}_4}(T)] \leq \sqrt{\frac{K}{r} T \log N} + rT = T^{2/3} (K \log N)^{1/3}$$

$$r^3 = \frac{K \log N}{T}$$

$$\frac{K T \log N}{r} \approx r^2 T^2 \Rightarrow r^3 \approx \frac{K}{T} \log N \quad rT = \frac{K^{1/3} (\log N)^{1/3}}{T^{1/3}}$$

Hedge bound $\sum_t E[\hat{c}_t^2(e_t) | P_t] \leq uT \approx \text{var of costs}$

$$r \propto \frac{1}{T^{1/3}} \quad rT \approx \frac{1}{T^{1/3}} \sqrt{KT \log N}$$

El K Hedge (1,1) ... (N, u, v, w) = ...

Better bound # $E[\hat{C}_t^2(e_t) | P_t] \leq \frac{K}{1-r} =: u$ allow for $r=0$.
 \Rightarrow skip uniform exp step in Exp(4)

For each arm a , let E_a be set of experts that recommend arm a .

$$\text{Let } p_t(a) = \sum_{e \in E_a} p_t(e)$$

$$q_t(a) \geq p_t(a) (1-r)$$

$$\Rightarrow \hat{C}_t(e) = \hat{C}_t(a_{t,e}) \leq \frac{c_t(a_{t,e})}{q_t(a_{t,e})} \leq \frac{1}{q_t(a_{t,e})} \leq \frac{1}{(1-r)p_t(a_{t,e})}$$

$$\begin{aligned} \sum_e \hat{C}_t^2(e_t) p_t(e) &= \sum_a \sum_{e \in E_a} p_t(e) \hat{C}_t(e) \hat{C}_t(e) \\ &\leq \sum_a \sum_{e \in E_a} \frac{p_t(e)}{(1-r)p_t(a_{t,e})} \hat{C}_t(a_{t,e}) \\ &= \frac{1}{1-r} \sum_a \frac{\hat{C}_t(a)}{p_t(a)} \underbrace{\sum_{e \in E_a} p_t(e)}_{p_t(a)} \\ &= \frac{1}{1-r} \sum_a \hat{C}_t(a) \end{aligned}$$

$$\begin{aligned} E[\hat{C}_t^2(e_t) | P_t] &= \frac{1}{1-r} \sum_a E[\hat{C}_t(a) | P_t] \\ &= \frac{1}{1-r} \sum_a c_t(a) \leq \frac{K}{1-r} \end{aligned}$$