

# Learning (and Sleeping) with Experts & Bandits

Online learning: ML framework based on sequential data (stocks, language, weather, ... most things) (has a certain "flavor" I want to communicate)

**PROBLEM** On day  $t=1, \dots, T$

- choose  $p_t \in \Delta_{N-1}$  ← prob simplex on  $[N]$  over "experts"
- sample  $i_t \sim p_t$   $i_t \in [N]$
- nature reveals  $l_t \in \mathbb{R}^N$
- learner incurs loss  $l_{t,i_t}$  (in expectation,  $\langle p_t, l_t \rangle$ )

crucial step →

Interpretation / setup: You have  $N$  experts  
 Every day, they make a prediction, and each has "loss"/error  $l_{t,i}$   
 Learner knows how experts have behaved in the past  
 maintains distribution  $p_t$

**GOAL** minimize  $\sum_{t=1}^T l_{t,i_t} - \min_{i \in [N]} \sum_{t=1}^T l_{t,i} = R_T$

in expectation,  $\mathbb{E}(R_T) = \sum_{t=1}^T \langle l_t, p_t \rangle - \min_{i \in [N]} \sum_{t=1}^T l_{t,i}$

**Algorithm: HEDGE** (MWU) 2

Maintain  $\eta > 0$ . Initialize  $w_1 = (1, \dots, 1)$   
 "weights"  
 On day  $t=1, \dots, T$ :

- choose  $p_t = \frac{w_t}{\|w_t\|_1} = \frac{w_t}{\sum w_{t,i}}$
- choose  $i_t \sim p_t$
- upon receiving  $l_t$

$$\forall i \in [N] \quad w_{t+1,i} = w_{t,i} \exp(-\eta l_{t,i})$$

Corollary

**Thm** Using HEDGE,  $\mathbb{E}(R_T) \leq \sqrt{2T \log N}$

with  $l_t \in [0,1]^N$

$$\sum_{t=1}^T \langle l_t, p_t \rangle - \min_{i \in [N]} \sum_{t=1}^T l_{t,i} \leq \frac{\ln N}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \langle l_t^2, p_t \rangle$$

Thm

(so average, per-round regret decays like  $\frac{1}{\sqrt{T}}$ )  $H(p) = \sum p_i \ln \frac{1}{p_i}$   
entropy of  $p$

**PF** observe that  $p_t = \min_{p \in \Delta^{N-1}} \left( \sum_{s=1}^{t-1} \langle l_s, p \rangle - \frac{H(p)}{\eta} \right)$

Standard proof: let  $Z_t = \sum_{i=1}^N w_{t,i}$  "potential function"

observe,  $\forall i, Z_{T+1} > w_{T+1,i} = \exp(-\eta \sum_{t=1}^T l_{t,i})$

$$\frac{Z_{t+1}}{Z_t} \leq \exp(-\eta \langle l_t, p_t \rangle + \frac{\eta^2}{2} \langle l_t^2, p_t \rangle)$$

$$\Rightarrow \ln(Z_{T+1}) - \ln(Z_1) \leq -\eta \sum_{t=1}^T \langle l_t, p_t \rangle + \frac{\eta^2}{2} \sum_{t=1}^T \langle l_t^2, p_t \rangle$$

link inequalities, rearrange  $\square$

**PROBLEM**

On day  $t=1, \dots, T$

- choose  $p_t \in \Delta_{N-1}$
- sample  $i_t \sim p_t$
- \* nature reveals ONLY  $l_{t,i_t}$
- incur loss  $l_{t,i_t}$

MINIMIZE  $\sum l_{t,i_t} - \min_{i \in [N]} \sum l_{t,i}$   
 SAME REGRET

3

"MULTI-ARMED BANDITS"

partial information

The only difference

ALGORITHM (explore + exploit) w/ exponential weights, EXP3

Maintain  $\eta > 0$ .  $w_1 = (1, \dots, 1)$

On day  $t=1, \dots, T$ :

- choose  $p_t = w_t / \|w_t\|_1$
- choose  $i_t \sim p_t$

$$w_{t,i} = \begin{cases} w_{t,i} & \text{if } i \text{ for } t \neq t \text{ for } i \neq i_t \\ w_{t,i} \exp(-\eta l_{t,i} / p_{t,i}) & \text{for } i = i_t \end{cases}$$

↑ inverse probability weighting

Idea: apply Hodge guarantee with loss vector

$$\tilde{l}_{t,i} = \mathbb{1}(i=i_t) l_{t,i} / p_{t,i}$$

$$\forall i \quad \sum_{t=1}^T \langle \tilde{l}_t, p_t \rangle \leq \sum_{t=1}^T \tilde{l}_{t,i} + \frac{\ln k}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \langle \tilde{l}_t^2, p_t \rangle$$

Take expectation:  $(\tilde{l}_{t,i} \rightsquigarrow l_{t,i} \quad \tilde{l}_t^2 \rightsquigarrow \tilde{l}_{t,i}^2 / p_{t,i})$

$$R_T \leq \frac{\ln k}{\eta} + \frac{\eta k}{2} \sum_{t=1}^T \underbrace{\mathbb{E} \langle \tilde{l}_t^2, p_t \rangle}_N \sum_{i=1}^N l_{t,i}^2$$

Note: We can decouple "action" and "expert"

for  $t=1, \dots, T$ : choose  $p_t \in \Delta_{k-1}$  ↓  $k$  actions ↓  $N$  experts

→ receive expert actions  $b_{t,i} \in [k]$  for  $i \in [N]$

learner chooses action  $a_t \sim p_t$

learner receives loss  $\ell_t(a_t) \in \mathbb{R}_+^k$

(in previous setup,  $k=N$  and no expert advice received)

Algorithm (EXP4):

for  $t=1, \dots, T$ :

- choose  $p_{t,a} = \frac{\sum_{i=1}^N w_{t,i} \mathbb{1}(b_{t,i}=a)}{\sum_{i=1}^N w_{t,i}} \in \Delta_{k-1}$

- choose  $a_t \sim p_t$ , receive  $\ell_{t,a_t}$

- create  $\tilde{\ell}_{t,a} = \mathbb{1}(a_t=a) \ell_{t,a_t} / p_{t,a} \quad \forall a \in [k]$

- update  $w_{t+1,i} = w_{t,i} \exp(-\eta \tilde{\ell}_{t,b_{t,i}})$   
 $i \in [N]$

Sleeping experts

a more insidious kind of partial information... [5]

On day  $t=1, \dots, T$

— choose  $p_t \in \Delta_{N-1}$

— nature reveals ONLY  $l_{t,i}$  for  $i \in E_t \subseteq [N]$

"awake experts"  
↓  
adversarially chosen

— incur loss  $\sum_{i \in E_t} p_{t,i} l_{t,i}$

Our approach: "sleepy hedge"

if  $i \notin E_t$ , keep weight unchanged

$i \in E_t$ , update  $w_{t+1,i} = w_{t,i} \exp(-\eta (l_{t,i} - \frac{1}{|E_t|} \sum_{j \in E_t} l_{t,j}))$

Guarantee:

$$R_T \leq \sqrt{2kNT \log N} + \sum_{t=1}^T \mu_t (p_{t,i^*} - \mathbb{1}(i^* \in E_t))$$

$(l_{t,i} - \mu_t)^2 \leq k \quad \forall t, i$

can "assume this away"