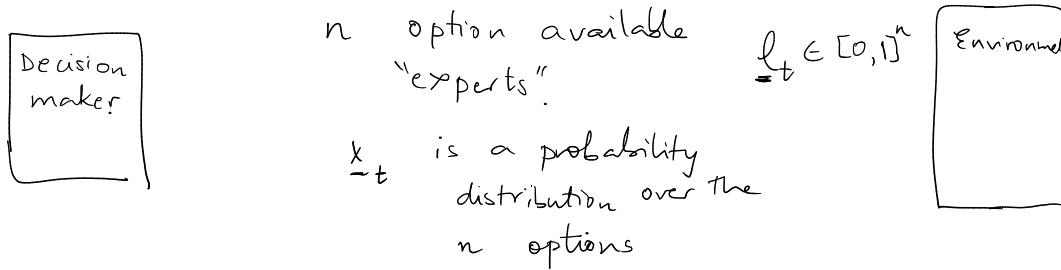


Online Learning with Expert Advice

(sequential decision making framework).

Multi-round game.



- At time time t , the loss incurred by the decision-maker is $\underline{x}_t \cdot \underline{l}_t$
- No assumption on the loss other than the fact that they are bounded.
- The loss vector \underline{l}_t is only revealed after the choice \underline{x}_t is made.
- Play this game for T rounds.

$$\text{Loss}(DM) = \sum_{t=1}^T \underline{x}_t \cdot \underline{l}_t$$

- "Best fixed strategy in hindsight" $\Delta_n = \left\{ \underline{x} \in \mathbb{R}^n \mid \underline{x} \geq 0 \text{ and } \underline{x} \cdot \underline{1} = 1 \right\}$
- $\underline{x}^* \in \arg \min_{\underline{x} \in \Delta_n} \sum_{t=1}^T \underline{x} \cdot \underline{l}_t = \arg \min_{\underline{x} \in \Delta_n} \underline{x} \cdot \underbrace{\sum_{t=1}^T \underline{x} \cdot \underline{l}_t}_{L_{T+1}}$
- Regret = $\text{Loss}(DM) - \min_{\underline{x} \in \Delta_n} \sum_{t=1}^T \underline{x} \cdot \underline{l}_t$

Goal : Minimise Regret.

- Regret grows sublinearly with T , $\text{Regret} = O(\sqrt{\log T})$.
- Regret can be negative

(Assume that the decision maker knows the time horizon T .)

(Take I): "Follow the leader".

At time t , pick $\underline{x}_t \in \arg\min_{\underline{x}} \sum_{s=1}^{t-1} \underline{x} \cdot \underline{l}_s = \arg\min_{\underline{x}} \underline{x} \cdot (\sum_{s=1}^{t-1} \underline{l}_s)$.

pick i s.t. $L_{t,i}$ is minimum, set $x_{t,i} = 1$, $x_{t,j} = 0$ $\forall j \neq i$.

Only have 2 options:

	ℓ_1	ℓ_2	\dots	\dots	ℓ_T
1	$\frac{1}{2}$	0	1	0	1
2	0	1	0	1	0

$$\text{Loss(DM)} = T, \quad \text{Loss(Best in hindsight)} = T/2, \quad \text{Regret} = T/2.$$

Algorithm: (Exponentially weighted Forecast / Multiplicative Weight Update / Hedge...)

- $\underline{\omega}_1 = (1, \dots, 1) \in \mathbb{R}^n$

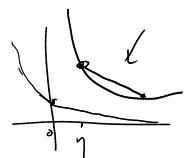
- For $t = 1, \dots, T$

- $\underline{x}_t = \frac{1}{Z_t} \underline{\omega}_t, \quad Z_t = \sum_i \omega_{t,i}$

- $\omega_{t+1,i} = \omega_{t,i} \cdot \exp(-\eta \ell_{t,i})$.

Theorem: Regret of MWUA with $\eta = \sqrt{\frac{2 \log n}{T}}$ is $2 \sqrt{(\log n) T}$.

Proof: $\frac{Z_{t+1}}{Z_t} = \frac{\sum_i \omega_{t+1,i}}{Z_t} = \frac{\sum_i \omega_{t,i} \exp(-\eta \ell_{t,i})}{Z_t} \leq \frac{\sum_i \omega_{t,i}}{Z_t} (1 + (e^{-\eta} - 1) \ell_{t,i}) = 1 + (e^{-\eta} - 1) \sum_i x_{t,i} \cdot \ell_{t,i}$



loss of decision maker at time t

$$\frac{Z_{t+1}}{Z_t} \leq \exp(x_t \cdot \ell_t (e^{-\eta} - 1)). \quad (\underline{A})$$

$$\frac{Z_{T+1}}{Z_1} = \prod_{t=1}^T \frac{Z_{t+1}}{Z_t} \leq \exp((e^{-\eta} - 1) \sum_{t=1}^T x_t \cdot \ell_t) \quad (\underline{A})$$

$$\forall i \quad \frac{Z_{T+1}}{Z_1} \geq \frac{\omega_{T+1,i}}{n} = \frac{\exp(-\eta \sum_{t=1}^T \ell_{t,i})}{n} \quad (\underline{B})$$

Let i^* be s.t. $\sum_{j=1}^T x_j^{i^*} = 1$, $x_j^{i^*} = 0$ for $j \neq i^*$ is in $\arg\min_x \sum_{t=1}^T l_t$.

$$\frac{\exp(-\eta \sum_{t=1}^T l_t, i^*)}{n} \leq \exp((e^{-1} - 1) \sum_{t=1}^T x_t \cdot l_t). \quad \text{Combining (A) \& (B)}$$

$$-\eta \sum_{t=1}^T l_t, i^* - \log n \leq (e^{-1} - 1) \sum_{t=1}^T x_t \cdot l_t$$

$$\eta \left(\underbrace{\sum_{t=1}^T x_t \cdot l_t}_{\text{Loss(DM)}} - \underbrace{\sum_{t=1}^T l_t, i^*}_{\text{Loss(Best in hindsight)}} \right) \leq (e^{-1} - (1-\eta)) \sum_{t=1}^T x_t \cdot l_t + \log n$$

$$\text{Regret} \leq \left(\frac{e^{-1} - (1-\eta)}{\eta} \right) \sum_{t=1}^T x_t \cdot l_t + \frac{\log n}{\eta}$$

$$\boxed{e^{-1} \leq 1 - \eta + \eta^2} \quad \leq \quad \eta^T + \frac{\log n}{\eta} \leq 2 \sqrt{T \log n} \quad \text{if } \eta = \sqrt{\frac{\log n}{T}}$$

Theorem: There exists a sequence of losses such that any algorithm suffers a regret of $\mathcal{O}(\sqrt{T})$. (or $\mathcal{O}(\sqrt{T \log n})$).

Proof ("sketch"):

At time t , toss a coin. if "HEADS" $l_t = (0, 1)$
else $l_t = (1, 0)$

$$\mathbb{E} [\text{loss of DM}] = T/2 \quad \mathbb{E} [\text{loss at time } t] = \frac{x_{t,1}}{2} + \frac{x_{t,2}}{2}$$

$$\mathbb{E} [\text{loss of best in hindsight}] = \frac{T}{2} - c\sqrt{T}$$

$$\mathbb{E} [\text{Regret}] = c\sqrt{T}$$

(Can prove that with some constant prob. $\text{Regret} \geq c\sqrt{T}$ for some constant c .)