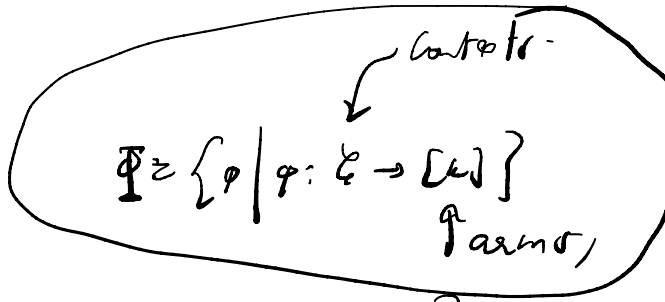




Contextual bandits:

Seeking to bound



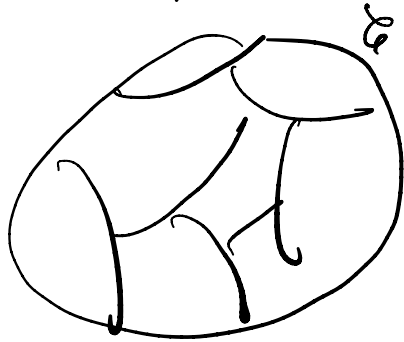
$$R_n = \mathbb{E} \left[ \max_{\phi \in \Phi} \sum_{t=1}^n x_t(\phi(\mathcal{C}_t)) - X_t \right]$$

(pseudo regret)

• Instead of playing with all functions from  $\mathcal{C} \rightarrow k$  we could restrict to a smaller subset of functions.

Regret will  $\downarrow$ ; or reward  $\uparrow$ ;

Examples:



$\mathcal{C}$ .

Partition  $\mathcal{C}$ ;

Play the same arm on all contexts in a given part;

Report will depend upon # parts;

(2) Similarity function:

$$s: \mathcal{C} \times \mathcal{C} \rightarrow [0, 1]$$

$$\bar{\Phi} = \left\{ \varphi: \mathcal{C} \rightarrow \mathcal{C} \mid \frac{1}{|\mathcal{C}|^2} \sum_{c, d \in \mathcal{C}} (1 - s(c, d)) \mathbb{1}[\varphi(c) \neq \varphi(d)] \leq \theta \in [0, 1] \right\}$$

- Dissimilarity function  $\in$  bounded;  
- Play against functions which will send similar contexts to the same arm.

- If similarity of contexts is low, then we will pull different arms for these contexts, but the # functions  $\varphi \in \bar{\Phi}$  is also not too many,  $\therefore 1 - s(c, d)$  is large;  $\therefore$  our report is better!

• Pick a collection of predictors,  $\phi_1, \dots, \phi_M$ ;  
each  $\phi_i: \mathcal{E} \rightarrow [k]^T$ ;

We can use a bandit algorithm to compete with the best of these in an online fashion;

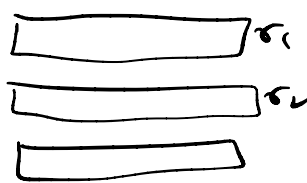
Adv:  $\phi_i$ 's can be obtained offline using batch training.

• So natural to analyze contextual bandits in the framework of experts;

VER ①: Simplest; Follow the expert.

k experts; sequence of rewards  $r_1, r_2, \dots$

$r_i \in [0, 1]^k$ ;



• Every day choose a distribution  $p_t$  over experts;

• At the end of the day see the quality of **EVERY** expert;

$$\text{Regret}(T) := \max_i \sum_{t=1}^T r_{t,i} - \sum_{t=1}^T \langle A_t, z_t \rangle$$

$$\text{or Regret}(T) := \max_P \sum_{t=1}^T \langle P, r_t \rangle - \sum_{t=1}^T \langle p_t, r_t \rangle$$

Best combination of experts;

Want a bound on min regret;

Start with a wt of 1 for each expert;  
 Given rewards,  $r_t^t$

$$w_i^{t+1} = w_i^t (1 + \epsilon r_i^t)$$

if  $r_i^t < \text{large}$   
 ↑ the wt of this expert;

- Given wts set  $\phi^t = \sum_i w_i^t$

• let  $\phi^t = \frac{\omega^t}{\phi^t}$ ,  $\omega^t = \langle \omega_{1,t}^t, \dots, \omega_{k,t}^t \rangle$

- we show that if the horizon is  $> \frac{\ln k}{\epsilon^2}$

then the regret  $\textcircled{2}$  satisfies,

$$\textcircled{2} \leq \frac{\ln k}{\epsilon} + \epsilon n \quad \leftarrow \text{horizon}$$

$$\frac{\textcircled{2}}{n} \leq \frac{\ln(k)}{n\epsilon} + \epsilon$$

average regret over the horizon  $\leq \epsilon + \frac{\ln(k)}{n\epsilon}$

set:  $n \geq \frac{\ln(k)}{\epsilon^2}$

AVG regret  $\leq \epsilon + \epsilon = 2\epsilon$

---

$$\phi^{t+1} = \sum w_i^{t+1} = \sum w_i^t (1 + \epsilon \delta_i^t)$$

$$= \sum p_i^t \phi^t (1 + \epsilon r_i^t)$$

$$\phi^{t+1} = \phi^t + \epsilon \phi^t \langle p^t, r^t \rangle$$

$$= \phi^t [1 + \epsilon \langle p^t, r^t \rangle]$$

$$\phi^{t+1} \leq \phi^t \cdot e^{\epsilon \langle p^t, r^t \rangle}$$

$$\leq \phi \cdot e^{\epsilon \sum \langle p_j^t, r_j^t \rangle}$$

$$= \phi \cdot e^{\epsilon \sum \langle p^t, r^t \rangle}$$

Now  $\phi^{t+1} \geq w_i^{t+1} \quad \forall i;$

$$= \prod_{t=1}^n (1 + \epsilon r_i^t);$$

close.

$$\epsilon \leq 1/2$$

$$\Rightarrow r_i^t \leq 1$$

$$\leq 1/2$$

Now for  $|x| \leq 1/2,$

$$\ln(1+x) \geq x - x^2$$

$$\phi^{t+1} \geq \prod_{t=1}^n e^{\epsilon r_i^t - \epsilon^2 (r_i^t)^2} = e^{\epsilon \sum r_i^t - \epsilon^2 \sum (r_i^t)^2}$$

$$k. \quad e^{\sum_{t=1}^n \langle p^t, r^t \rangle} \geq e^{\sum_{t=1}^n r_i^t} \cdot e^{-\sum_{t=1}^n (r_i^t)^2}$$

$$\ln k + \sum_{t=1}^n \langle p^t, r^t \rangle \geq \sum_{t=1}^n r_i^t - \sum_{t=1}^n (r_i^t)^2$$

$$\frac{\ln(k)}{\varepsilon} + \sum_{t=1}^n \langle p^t, r^t \rangle \geq \sum_{t=1}^n r_i^t - \varepsilon \sum_{t=1}^n (r_i^t)^2$$

or:  $\sum_{t=1}^n r_i^t - \sum_{t=1}^n \langle p^t, r^t \rangle \leq \frac{\ln(k)}{\varepsilon} + \varepsilon \sum_{t=1}^n (r_i^t)^2$

$$\sum_{t=1}^n r_i^t - \sum_{t=1}^n \langle p^t, r^t \rangle \leq \frac{\ln(k)}{\varepsilon} + \varepsilon n \quad \uparrow \quad [0, 1]$$

Taking a convex comb of these inequalities

$$\max_p \sum_{t=1}^n \langle p, r^t \rangle - \sum_{t=1}^n \langle p^t, r^t \rangle \leq \frac{\ln(k)}{\varepsilon} + \varepsilon n$$

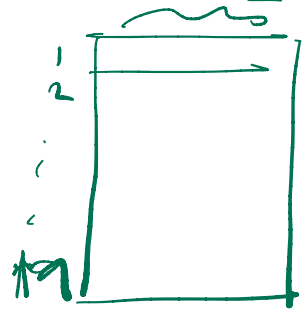
$$\text{Avg} \leq \frac{\ln(k)}{\varepsilon} + \varepsilon \quad \underline{\text{QED}}$$



# Repeat analysis with expert:

- 1) Input  $n, k, M, \eta, \delta$
- 2)  $Q_i = (1/M, \dots, 1/M)$
- 3) for  $t=1$  to  $n$
- 4) Get advice  $E^{(t)}$
- 5) Choose action  $A_t \sim P_t, P_t = Q_t E^{(t)}$
- 6) Get  $X_t = x_{tA_t}$
- 7) Set  $\hat{X}_{t,i} = 1 - \frac{\mathbb{1}(A_t=i)(1-x_t)}{P_{t,i} + \delta}$
- 8) Propagate rewards  $\tilde{X}_{t,i} = E^{(t)A_t} X_{t,i}$
- 9)  $Q_{t+1,i} = \frac{\exp(\eta \tilde{X}_{t,i}) Q_{t,i}}{\sum_j \exp(\eta \tilde{X}_{t,j}) Q_{t,j}} \forall i \in [k]$
- 10) end for

$M = \# \text{ experts.}$



• Analyze when  $\delta = 0$ :

Then:  $\delta = 0$ ;  $\eta = \sqrt{2 \log(M) / nk}$  & let

$R_n = \text{ regret after } n \text{ rounds}$ ; Then

$$R_n \leq \sqrt{2nk \log(M)} \iff \sqrt{\underbrace{nk}_{=T} \log \underbrace{k}_{\uparrow} \underbrace{M}_{\downarrow}}$$

Recall we showed:  $\forall i$

$$\hat{s}_{ni} - s_n \leq \frac{\log(k)}{\eta} + \eta \sum_{t=1}^n \sum_{j=1}^k P_{tj} \hat{x}_{tj}^2$$

that proof yields:

stronger to  $\frac{\eta}{2} \sum_{t=1}^n \sum_{j=1}^k P_{tj} \hat{y}_{tj}^2$

$$\sum_{t=1}^n \underbrace{\hat{x}_{t,m}^2}_{\sim} - \sum_{t=1}^n \sum_{m=1}^M \underbrace{P_{tm} \hat{x}_{tm}^2}_{\sim} \leq \frac{\log(M)}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \sum_{m=1}^M P_{tm} (1 - \hat{x}_{tm}^2)$$

$E^{(1)}, A_1, \dots, A_{t-1}, E^{(t)}$  be the history,

$$\mathbb{E}_t[\ ] = \mathbb{E}[\ ] \mid \text{history}$$

Let  $m^* = \underset{m \in [M]}{\operatorname{argmax}} \sum_{t=1}^n \mathbb{E}_m^t x_t$  ↖ expected return of expert  $m$ .

Now:  $\mathbb{E} = \mathbb{E}[\mathbb{E}_t]$

$$\left. \begin{aligned} \sum_{t=1}^n \underbrace{\tilde{x}_{t,m^*}} - \sum_{t=1}^n \sum_{m=1}^M p_{t,m} \tilde{x}_{t,m} &\leq \frac{\log(M)}{\eta} \\ + \frac{\eta}{2} \sum_{t=1}^n \sum_{m=1}^M p_{t,m} (1 - \tilde{x}_{t,m})^2 & \end{aligned} \right\} \quad \color{red}{***}$$

$\tilde{x}_{t,i}$  - unbiased ( $\sigma = 0$ )  $\therefore \mathbb{E}_t[\tilde{x}_t] = x_t$

$$\underbrace{\mathbb{E}_t[\tilde{x}_t]} = \mathbb{E}_t[\mathbb{E}^{(t)}[\tilde{x}_t]] = \mathbb{E}^{(t)}[\underbrace{\mathbb{E}_t[\tilde{x}_t]}_{x_t}]$$

Now

$$\mathbb{E}[\tilde{x}] = \mathbb{E}[\mathbb{E}_t[\tilde{x}]]$$

Applying  $\mathbb{E}$  on both sides and using the above, (first apply  $\mathbb{E}_t$  & then  $\mathbb{E}$  on both sides)

$$\mathbb{E} \left[ \mathbb{E}^{(t)}_{\alpha_t} - \sum_{t=1}^n \underline{x}_t \right] \stackrel{= R_n}{\leq} \frac{\log M}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \sum_{m=1}^M \mathbb{E} \left[ \mathbb{Q}_{t,m} (1 - \tilde{x}_{t,m}) \right]$$

$$\therefore R_n \leq \frac{\log M}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \sum_{m=1}^M \mathbb{E} \left[ \mathbb{Q}_{t,m} (1 - \tilde{x}_{t,m}) \right]$$

Work with  $\hat{y}_{t,i} = 1 - \hat{x}_{t,i}$ ;  $y_{t,i} = 1 - \alpha_{t,i}$

$$\tilde{y}_{t,m} = 1 - \tilde{x}_{t,m}$$

Now:  $\tilde{y}_t = \mathbb{E}^{(t)} \hat{y}_t$

Use notation  $A_{t,i} \stackrel{\Delta}{=} \mathbb{1}\{A_t = i\}$

Then  $\hat{y}_{t,i} = \frac{A_{t,i} y_{t,i}}{p_{t,i}}$

$$\mathbb{E}_t \left[ \sum_{t=1}^k \tilde{y}_{tm}^2 \right] = \mathbb{E}_t \left[ \left( \frac{\sum_{i=1}^k E_{mi}^{(t)} y_{tAT}}{P_{tAT}} \right)^2 \right]$$

$$= \sum_{i=1}^k \frac{(E_{mi}^{(t)} y_{tAT})^2}{P_{ti}} \leq \sum_{i=1}^k \frac{E_{mi}^{(t)}}{P_{ti}}$$

$$\therefore \mathbb{E} \left[ \sum_{m=1}^M \mathcal{Q}_{tm} (1 - \tilde{x}_{tm})^2 \right]$$

$$= \mathbb{E} \left[ \mathbb{E}_t \left[ \sum_{m=1}^m \mathcal{Q}_{tm} (\tilde{y}_{tm}^2) \right] \right]$$

$$\leq \mathbb{E} \left[ \sum_{m=1}^m \mathcal{Q}_{tm} \sum_{i=1}^k \frac{E_{mi}^{(t)}}{P_{ti}} \right]$$

$$= \mathbb{E} \left[ \underbrace{\sum_{i=1}^k \frac{\sum_{m=1}^m \mathcal{Q}_{tm} E_{mi}^{(t)}}{P_{ti}}}_{(1)} \right] = \underline{k} \quad \because P_t = \mathcal{Q}_t E^{(t)}$$

$$R_n \leq \frac{\log M}{\eta} + \frac{\eta n k}{2}$$

$$\eta = \sqrt{\frac{2 \log M}{nk}}$$

$$= \underline{\underline{\sqrt{2 \eta n k \log M}}}$$