

## 0-sum games and MW

Recall: given a payoff matrix  $M$  these are payoffs from column player  $C$  to row player  $R$ .

$(M)_{m \times n}$ . If row player plays  $x$  and column player plays  $y$  then

$$E[\text{payoff}] = x^T M y. \quad x \in \Delta_m, y \in \Delta_n.$$

For a fixed play  $x$  by row player, let  $C(x)$  be the best response by col. player.

$$C(x) = \min_{y \in \Delta_n} x^T M y = \min_{e_j} x^T M e_j \leftarrow \text{since this is a vertex of the probability simplex.}$$

Similarly, let

$$R(y) = \max_{x \in \Delta_m} x^T M y = \max_{e_j} e_j^T M y.$$

Thm [Von Neumann]:  $\max_{x \in \Delta} C(x) = \min_{y \in \Delta} R(y).$   $[x, y, C(x) \leq R(y)]$

For any finite matrix  $M$ . Call this the value of the game  $V$ .

The online learning guarantee can prove this for us.

This is also provable via strong LP duality [HW].

Let's look at a gain version of the online learning problem: - gain vectors  $g^{(k)} \in [-1, 1]^N$ .

$\forall \epsilon, \forall T \geq \frac{4gN}{\epsilon^2}, \forall i$

$$\frac{1}{T} \sum_t \langle g^t, p^t \rangle \geq \frac{1}{T} \sum_t \langle g^t, e_i \rangle - \epsilon$$

Nb. deterministic, so can assume adaptive  $g^t$  based on  $p^t$ .

assuming gain vectors are in  $[-1, 1]$ .

So now the algorithm is:-

Assume payoffs are in  $[-1, 1]$ .

- the  $m$  pure strategies of row player are experts. Start with uniform ~~static~~ weights.

- Each time the column player plays best response ~~to  $p^t$~~ . say column  $j(t)$   
 $\rightarrow e_{j(t)}$  to  $p^t$ . Set  $g^t = M e_{j(t)}$ .

After  $T \geq \Omega\left(\frac{\log m}{\epsilon^2}\right)$  steps, let  $\hat{z} = \frac{1}{T} \sum p^t$  and  $\hat{y} = \frac{1}{T} \sum g^t$ .

Claim:  $C(\hat{x}) \leq R(\hat{y}) \leq C(\hat{x}) + \epsilon$  for  $T \geq \frac{\Omega(\log m)}{\epsilon^2}$ .

Pf: the first inequality is true for all  $x, y \in \Delta_n$

[i.e. the  $C(\hat{x})$  is what the row player can guarantee even by playing first with  $\hat{x}$

$R(\hat{y})$  is what the col player can enforce even by playing first with  $\hat{y}$ ].

Other direction: Recall  $\hat{x} = \frac{1}{T} \sum_t p^t$   $\hat{y} = \frac{1}{T} \sum_t q^t$ .

At each time payoff is  $(p^t)^T M q^t = C(p^t) \stackrel{\text{by best response of col}}{\leq} C(\hat{x})$

$$\Rightarrow \frac{1}{T} \sum_t (p^t)^T M q^t \leq \frac{1}{T} \sum_t C(p^t) \leq C\left(\frac{\sum_t p^t}{T}\right) = C(\hat{x})$$

but now use the gains version of online regret bound

$$\Rightarrow \frac{1}{T} \sum_t \langle e_i, M q^t \rangle \leq C(\hat{x}) + \epsilon \quad \forall i$$

$$\Rightarrow \sum_t \langle e_i, M \hat{y} \rangle \leq C(\hat{x}) + \epsilon \quad \forall i$$

$$\Rightarrow R(\hat{y}) \leq C(\hat{x}) + \epsilon. \quad \square$$

So: one side does online learning, the other side plays best response, then we get ~~max~~ minimax theorem, using compactness.

(since for any  $\epsilon$ , we can ensure that  $C(x)$  and  $R(y)$  are closer than  $\epsilon$ ).