

Reminder: Presentations every Th

(Swap Steve in April
aiming)

Assignment: Volunteer by Wed

Bobby + Dan midnight or get randomly assigned.

1st 2 weeks: Kalai-Vempala framework
(Exp. many experts)

Last 2 weeks: Multi-armed bandit

Online learning & games

low regret algo \rightarrow converges to NE
for 2 player 0-sum
games.

Q: _____ \rightarrow general non 0-sum?

(Not known)

Notes: heuristic argument against this.

Q: low regret \rightarrow some stable configuration

Blackwell's approachability thm

low "internal" regret \rightarrow Correlated Equilibria
(Aumann?)

$$g = (k, A_i, u_i)$$

$$NE: p \in \prod_{i=1}^k \Delta(A_i) \quad (\text{product distribution})$$

$$CE: p \in \Delta\left(\prod_{i=1}^k A_i\right)$$

p is a CE: Trusted 3rd party
3rd party will sample a pure strategy according to p → no incentive for player to deviate

Defn: $P \in \Delta \left(\prod_{i=1}^k A_i \right)$ is a CE

if $\forall i \in [k], \forall \alpha, \beta \in A_i$ s.t. $\Pr[i \text{ picks } P \cdot \alpha] > 0$

$$\mathbb{E}_P [u_i(\alpha, a_{-i}) \mid a_i = \alpha]$$

$$\geq \mathbb{E}_P [u_i(\beta, a_{-i}) \mid a_i = \alpha]$$

\equiv
 $\sum_{a_{-i}}$

$$(u_i(\alpha, a_{-i}) - u_i(\beta, a_{-i})) \cdot \Pr_P[\alpha, a_{-i}] \geq 0$$

$$\left[\sum_{a \in \prod_i A_i} \Pr[a] = 1 \right] \forall a \Pr[a] \geq 0$$

Player 1

		Player 2	
		Stop	Go
Stop	(4, 4)	(1, 5)	
Go	(5, 1)	(0, 0)	

All NE:

0	1
0	0

0	0
1	0

.25	.25
.25	.25

Some CE:

0	.5
.5	0

$\frac{1}{3}$	$\frac{1}{3}$
$\frac{1}{3}$	0

(External) Regret : Loss compared to
a constant sequence.

(Internal) Regret : Loss compared to a
"class" of sequences

External regret \leq Internal regret

Internal Regret

$A \rightarrow$ set of actions & given any sequence

$\bar{a} = a_1, \dots, a_T \in A$ & sequence of

payoff functions $g_1, \dots, g_T: A \rightarrow \mathbb{R}$

& an "advisor" function $f: A \rightarrow A$

$$\hat{R}_f(\bar{a}, \bar{g}, T) = \frac{1}{T} \sum_{t=1}^T g_t(\bar{a}_t) - g_t(f(a_t))$$

$$\hat{R}_{\text{int}}(\bar{a}, \bar{g}, T) = \max_{f: A \rightarrow A} \hat{R}_f(\bar{a}, \bar{g}, T)$$

$$\hat{R}_{\text{ext}}(\bar{a}, \bar{g}, T) = \max_{b \in A} \hat{R}_b(\bar{a}, \bar{g}, T)$$

$$\Rightarrow \hat{R}_{\text{ext}}(\bar{a}, \bar{g}, T) \leq \hat{R}_{\text{int}}(\bar{a}, \bar{g}, T)$$

Pairwise regret:

$$\hat{R}_{a,b}(\bar{a}, \bar{g}, T) = \frac{1}{T} \sum_{t=1}^T (g_t(b) - g_t(i_t)) \mathbb{1}_{i_t=a}$$

$$\hat{R}_{\text{pair}}(\bar{a}, \bar{g}, T) = \max_{a,b} \hat{R}_{a,b}(\bar{a}, \bar{g}, T)$$

Lemma: $\hat{R}_{\text{pair}} \leq \hat{R}_{\text{int}} \leq |A| \cdot \hat{R}_{\text{pair}}$

Def: \bar{a} has no internal regret
wrt \bar{g} if $\lim_{T \rightarrow \infty} \hat{R}_{\text{int}}(\bar{a}, \bar{g}, T) \rightarrow 0$

← Lemma implies

no internal regret

iff $\lim_{T \rightarrow \infty} \hat{R}_{\text{pair}}(a, \bar{g}, T) \rightarrow 0$

$$\hat{R}_{\text{pair}}(\bar{a}, \bar{g}, T) \leq \hat{R}_{\text{int}}(\bar{a}, \bar{g}, T) \leq |A| \cdot \hat{R}_{\text{pair}}$$

Pf: Fix $f: A \rightarrow A$

$$\hat{R}_f(\bar{a}, \bar{g}, T) = \frac{1}{T} \sum_{t=1}^T (g_t(f(a_t)) - g_t(a_t))$$

$$= \frac{1}{T} \sum_{t=1}^T \sum_{a \in A} (g_t(f(a)) - g_t(a)) \cdot \mathbb{1}_{a_t=a}$$

$$= \sum_{a \in A} \frac{1}{T} \sum_{t=1}^T (g_t(f(a)) - g_t(a)) \cdot \mathbb{1}_{a_t=a}$$

$$= \sum_{a \in A} \hat{R}_{g, f(a)}(\bar{a}, \bar{g}, T) \leq |A| \cdot \hat{R}_{\text{pair}}^x(\bar{a}, \bar{g}, T)$$

No internal regret \implies CE

Def: $S \subseteq \mathbb{R}^n$, $x \in \mathbb{R}^n$

$$\text{dist}(x, S) = \min_{y \in S} \|x - y\|_2$$

x_1, \dots, x_n converges to S if $\lim_{n \rightarrow \infty} \text{dist}(x_n, S) = 0$

Let $k, |A_i| < \infty$ $g = (k, A_i, u_i)$

Suppose each i repeatedly plays g using a no-internal regret learning algorithm with $g_t^i(a) = u_i(a, a_{-i}^t)$ to generate the sequence a_i^t for $t=1, 2, \dots$

Let $\bar{p}(T)$ be the uniform distribution over the multiset $\{(a_1^t, \dots, a_k^t) \mid 1 \leq t \leq T\}$

Then $\bar{p}(T)$ converges to C , the set of CE of g .

(E: $\mathbb{P} \in \Delta(\prod_{i \in [k]} A_i)$ is a CE

$\iff \forall i \in [k]; \alpha, b \in A_i$ s.t.

$$\sum_{\substack{a_{-i} \\ \in \prod_{j \neq i} A_j}} (u_i(b, a_{-i}) - u_i(\alpha, a_{-i})) \cdot \Pr[\alpha, a_{-i}] \cdot \Pr \text{ of } i \text{ playing} \geq 0$$

≤ 0

$a = (\alpha, a_{-i})$

For contradiction assume $\bar{p}(T)$ does not converge to \mathcal{C} .

I.e. let p be the limit of the sequence $\bar{p}(T_1), \bar{p}(T_2), \dots$

then $\text{dist}(p, \mathcal{C}) \geq \delta > 0$.

$\Rightarrow p \notin \mathcal{C}$

$\Rightarrow \exists i \in [k], \alpha, \beta \in A_i$ s.t.

$$\sum_{a-i} \left[u_i(\beta, a-i) - u_i(\alpha, a-i) \right] p(\alpha, a-i) = \varepsilon > 0.$$

As p is the limit pt of $\bar{p}(T), \dots$

large
 $\exists s \geq 1$ s.t.
 enough

$$\sum_{a-i} [u_i(\beta, a-i) - u_i(\alpha, a-i)] \bar{p}(T_s)$$

$$\Rightarrow \frac{1}{T_s} \sum_{a-i} \left([u_i(\beta, a-i) - u_i(\alpha, a-i)] \sum_{t=1}^{T_s} \mathbb{1}_{a^t=a} \right) \geq \frac{\epsilon}{2} > 0.$$

$$\mathbb{1}_{a^t=a} = \mathbb{1}_{a_i^t=\alpha} \cdot \mathbb{1}_{a_{-i}^t=a-i} \quad (a^t = (a_1^t, \dots, a_{i-1}^t, a_i^t, \dots, a_K^t))$$

$$\frac{1}{T_s} \sum_{t=1}^{T_s} \left(\sum_{a-i} [u_i(\beta, a-i) - u_i(\alpha, a-i)] \right)$$

$$\Rightarrow \frac{1}{T_s} \sum_{t=1}^{T_s} (u_i(\beta, a^t-i) - u_i(\alpha, a^t-i)) \geq \frac{\varepsilon}{2}$$

$\mathbb{1}_{a^t-i = a-i} \cdot \mathbb{1}_{a^t-i = \alpha}$

$\mathbb{1}_{a^t-i = \alpha} \geq \frac{\varepsilon}{2}$

$$\hat{P}_{\alpha, \beta} \left((a_i^1, a_i^2, \dots, a_i^{T_s}), (g_i^1, \dots, g_i^{T_s}) \right)$$

Contradiction! $\geq \frac{\varepsilon}{2} > 0$.