

# 6.S890: Topics in Multi-Agent Learning

## Lecture 9: Stochastic Games Supplement

↳ Existence of Nash Equilibrium in infinite horizon, discounted, stochastic games w/ finite states and actions

Theorem: Every infinite-horizon discounted stochastic game w/ a finite # players, # states, # actions has a Nash equilibrium in stationary Markov policies.

That is there  $\exists$  a collection of policies  $\pi_1, \dots, \pi_m$  where  $\pi_i: S \rightarrow \Delta(A_i)$ ,  
such that:  $\uparrow$  states  $\uparrow$  actions of player  $i$

$$u_i(\pi_i; \pi_{-i}) \geq u_i(\pi'_i; \pi_{-i}), \forall i, \forall \pi'_i$$

↳ even if allowed to be history-dependent

Proof: For policy profile  $\pi = (\pi_1, \pi_2, \dots, \pi_m)$   
player  $i \in \{1, \dots, n\}$

we'll define  $(v_i^\pi(s))_{s \in S}$

and  $(q_i^\pi(s, a_i))_{s \in S, a_i \in A_i}$

as follows:

- $v_i^\pi(s)$ : infinite discounted utility of player  $i$  if game started at state  $s$  & players used policies  $\pi_1, \dots, \pi_m$

$$\hookrightarrow v_i^\pi(s) = \underbrace{\sum_a r_i(s, a) \cdot \pi(a|s)}_{\stackrel{\text{IID}}{r_i^\pi(s)}} + \sum_{s'} v_i^\pi(s') \cdot \underbrace{\gamma \cdot \sum_a \pi(a|s) \cdot P(s'|s, a)}_{\stackrel{\text{IID}}{\Gamma^\pi(s, s')}} \quad \leftarrow \text{discount factor}$$

$$\Rightarrow (\mathbf{I} - \gamma \cdot \Gamma^\pi) v_i^\pi = r_i^\pi \quad \Rightarrow v_i^\pi = (\mathbf{I} - \gamma \cdot \Gamma^\pi)^{-1} \cdot r_i^\pi$$

⊛ why is  $\mathbf{I} - \gamma \cdot \Gamma^\pi$  invertible?

$$\text{Row sums of } \Gamma^\pi: \sum_{s'} \Gamma^\pi(s, s') = \sum_a \pi(a|s) = 1$$

$\gamma < 1$   
 $\Rightarrow \mathbf{I} - \gamma \cdot \Gamma^\pi$  is strictly diagonally dominant

$\Rightarrow \mathbf{I} - \gamma \cdot \Gamma^\pi$  is non-singular

by above note that, as a f'n of  $\pi$ ,  $v_i^\pi(s)$  is continuous

- $q_i^\pi(s, a_i)$ : infinite discounted utility of player  $i$  if game started at  $s$ , players used policies  $\pi_1, \dots, \pi_m$  throughout, except at the very first step player  $i$  plays  $a_i$

$$q_i^\pi(s, a_i) = \sum_{a_{-i}} r_i(s, a) \cdot \pi_{-i}(a_{-i} | s) + \sum_{s'} v_i^\pi(s') \cdot \gamma \cdot \sum_{a_{-i}} \pi_{-i}(a_{-i} | s) \cdot P(s' | s, a)$$

↳ as a function of  $\pi$   $q_i^\pi(s, a_i)$  is continuous b.c. everything on the RHS is cont. wrt.  $\pi$  including  $v_i^\pi(s)$

- Now define Nash-type function mapping  $\pi \xrightarrow{f} \pi'$  as follows:

$$\forall i, \forall s, \forall a_i:$$

$$\pi'_i(a_i | s) \leftarrow \frac{\pi_i(a_i | s) + \max(0, q_i^\pi(s, a_i) - v_i^\pi(s))}{1 + \sum_{a'_i} \max(0, q_i^\pi(s, a'_i) - v_i^\pi(s))}$$

$f$ : continuous over convex, compact set  $\times \underbrace{(\Delta(A_i))_i^S}$

$\Rightarrow \exists$  fixed point  $\pi = f(\pi)$   
Brouwer

↖ set of all possible stationary Markov policies of player  $i$

Claim:  $\pi$  is Nash Equilibrium

Proof: suffices to prove  $\pi_i$  is best response to  $\pi_{-i}$  among all stationary Markovian policies

[why? b.c. fixing  $\pi_{-i}$ , player  $i$  faces Markov Decision process and it's known that MDPs have optimal policies that are stationary & Markovian]

Interesting to note: 
$$V_i^\pi(s) = \sum_{a_i} \pi_i(a_i) Q_i^\pi(s, a_i) \quad (1)$$

↳ thus by doing the same logic we did in Nash's proof for each  $s$  separately it follows from  $\pi = f(\pi)$  that:

$$\forall s, \forall i, \forall a_i: Q_i^\pi(s, a_i) = 0$$

$$\Leftrightarrow \forall s, \forall i, \forall a_i: V_i^\pi(s) \geq Q_i^\pi(s, a_i)$$

⊛⊛

⊛⊛  $\Leftrightarrow \pi$  is a Nash Equilibrium.

$$\text{⊛⊛} \quad \forall i, \forall s, \forall a_i: \underbrace{V_i^\pi(s)}_{\text{player } i\text{'s expected infinite discounted payoff from } \pi_i} \geq \underbrace{Q_i^\pi(s, a_i)}_{\text{player } i\text{'s expected infinite discounted payoff at } s \text{ from single-round deviation from } \pi_i} \quad (2)$$

- clearly implied by  $\pi$  being a Nash Eq (so  $\Leftarrow$  easy)
- why does it imply Nash eq?

fixing  $\pi_{-i}$ , player  $i$  faces an MDP where:

$$\text{MDP} \quad \tilde{r}_i(s, a_i) = \sum_{a_{-i}} r(s, a) \cdot \pi_{-i}(a_{-i} | s)$$

$$\tilde{P}(s' | s, a_i) = \sum_{a_{-i}} P(s' | s, a) \pi_{-i}(a_{-i} | s)$$

given a policy  $\pi_i$ , denote by  $\tilde{V}_i^{\pi_i}(s)$  the expected infinite discounted payoff starting at  $s$  & using  $\pi_i$  in  $\tilde{\text{MDP}}$

(2) + (1)  $\Rightarrow \pi_i$  satisfies what are called "Bellman equations"

namely:

$$\forall s: \tilde{v}_i^{\pi_i}(s) = \max_{a_i} \left( \tilde{r}_i(s, a_i) + \gamma \cdot \sum_{s'} \tilde{P}(s'|s, a_i) \tilde{v}_i^{\pi_i}(s') \right)$$

Bellman Equations  $\Rightarrow \pi_i$  is optimal in  $\tilde{M}DP$   
so  $\pi_i$  is best response to  $\pi_i$ .

⊠