

# Problem Setting ; ST Game model - MAREL Paper

$$\mathcal{M} := \{S, \{A_i\}_{i=1}^n, P, r, \rho, H\}$$

→  $n$  agents

→ ' $H$ ' → horizon of Markov Game

$P$   
↓  
defines state  
transition prob function,  
defines game state's  
change over time

→  $s_h \in S$  → all system states

→ state at horizon

→  $A_i$  (action space of agent  $i$ )

$A = A_1 \times A_2 \times \dots \times A_n$  (combination of entire  
action space of all  
agents)

→ going to  $S_{h+1}$  from  $S_h$  at time step ' $h+1$ '

A actions  $a_h = \{a_{1,h}, a_{2,h}, \dots, a_{n,h}\}$  taken by all ' $n$ '  
agents at time ' $h$ '

next state,  $s_{h+1} \sim P(\cdot | s_h, a_h)$ , drawn  
from prob distribution defined by ' $P$ ',  
given the current state /  $s_h$  & actions  $a_h$

' $r$ ' is reward & ' $r_h$ ' is reward at time step ' $h$ '

$$r_h : S \times A \rightarrow [0, 1] \text{ matlab ke } r_h(s_h, a_h)$$

' $P$ ' → initial state dist

is reward  
that we get at a

Stochastic policy  $\pi = \{\pi_h : S \rightarrow \Delta(A)\}_{h=1}^H$

specific ' $h$ ', bounded  
b/w 0 & 1.

→ is a strat for agents which  
specifies how to choose actions

probabilistically at each stage in the game

$$\Delta(A)$$

→ prob simplex over  
action space ' $A$ ' → all possible prob  
dists to distribute  
the action space  $A$   
over  
means that  
agents choose their  
actions non-deterministically  
but by probabilities defined  
by this policy.



## Understanding the Value function & Q-function

Value function.  $V_{i,h}^{\pi}(s) \Rightarrow$  means value function for  
 i<sup>th</sup> player, at time step  $h$ ,  
 & state  $s$  and under the  
 ST policy  $\pi$  defines  
 The expected cumulative reward  
 that agent 'i' receives, defined

as,

$$V_{i,h}^{\pi}(s) := E \left[ \sum_{h'=h}^H r_{i,h'}(s_{h'}, a_{h'}) \mid s_h = s \right]$$

$\nwarrow$  expectation  
 $\swarrow$  cumulative reward over time from  $h'$  till total time  $H$

Q-function for agent 'i' at  
 time step  $h$  & state  $s$  and so on...

basically represents the expected  
 cumulative reward that agent 'i'  
 receives starting from state 's',

defined as

$$Q_{i,h}^{\pi}(s, a) = E \left[ \sum_{h'=h}^H r_{i,h'}(s_{h'}, a_{h'}) \mid s_h = a, a_h = a \right]$$

$\nwarrow$  depends on  $\pi_{h:H}$   
 $\swarrow$  This expectation is taken over future rewards, assuming that policy ' $\pi$ ' gets followed for all future stages

MPE defined  
 → deterministic strategy  
 Markov perfect Equilibrium

$$Q_{i,h}^{\pi}(s, a^*) > Q_{i,h}^{\pi}(s, a_i, a_{-i}^*) \quad \forall a_i \neq a_i^*$$

(basically means that Q-value (or expected reward) of choosing equilibrium action  $a^*$  must

be strictly greater than Q-value of any alt

→ essentially ensures that no agent deviates from their equilibrium action as all other alt actions yield lower rewards.



→ Definition 2 [✓]

→ Definition 3 [✓]

→ Definition 4 [✓]

→ 2.2 [Equilibrium Selection for Normal Form Games]

$r_i: A \rightarrow \mathbb{R}$ : the reward fine for agent  $i$  (total 'n' agents), depends on the joint action space  $A = A_1 \times A_2 \times \dots \times A_n$

→ Iterative Learning,

→ these agents adjust strats iteratively based on past st outcomes  
→ learning rules adjusted by Markov Chain

Ex 1  
'LL'

→ No aux vars,  $\{E = \emptyset\}$

→ transition probs based on agent 'i' rewards

$$K_{\epsilon}^{(t+1)} = (a_{-i}^{(t+1)}, a_i^{(t+1)}) | a_i^{(t)}$$

→ basically actions with higher rewards are more probable/likely but not guaranteed.

→ as  $\epsilon \rightarrow 0$ , we get to best response strat

→ as  $\epsilon > 0$ , learning process is 'ergodic', ensuring convergence to a unique stationary dist.

$$(a^{(t+1)}, \epsilon^{(t+1)}) \sim K_{\epsilon}(i, \cdot | a^{(t)}, \epsilon^{(t)})$$

the main transition kernel  
action taken at iteration 't'  
auxiliary vars assisting in learning

→ defines the prob of transitioning to next state given current state.

→ 'ε' is the 'rate of mistakes', adds randomness in learning

$$\{r_i\}_{i=1}^n = e^{-r_i(a_i^{(t+1)})} / \sum_{a_i} e^{-r_i(a_i, a_{-i}^{(t)})}$$

→ main purpose of iterative learning

introduction of 'ε' or rate of mistakes results in agents converging to a 'good' equilibrium

→ Assumption 1 (Ergodicity) [✓]

→ predictable sys cuz as  $t \rightarrow \infty$ , prob of being in any state converges to a fixed value regardless of initial state.

guarantees stable & unbiased equilibrium selection.

→ system is a learning agent navigating a map of possible decisions, a lot of exploration element (agent explores & can also make mistakes). Over time, agent learns

& and develops an understanding of best places to stay

Definition 5 [✓] / Definition 6 [✓]

SSE → stochastically stable equilibria → more stable NEs

' $\epsilon$ ' → hidden var defined in Ex 2

→ Assumption 2

resistance ' $R$ ' is a constant such that, along with constants  $C_1, C_2$  ( $C_1, C_2 > 0$ )

→ measures cost of transitioning b/w states in learning process

$$R((a, \xi) \rightarrow (a', \xi'))$$

↪ linked to transition probability

$$C_1 \epsilon^{R((a, \xi) \rightarrow (a', \xi'))} < K^\epsilon(a', \xi' / a, \xi; \{\pi_i\}_{i=1}^n)$$

→ higher the value of ' $R$ ,'

harder the transition b/w states

→  $R=0$ , transition similar to when  $\epsilon \rightarrow 0$

→  $R=\infty$ , transition is as impossible as when  $K_\epsilon = 0$  for all  $\epsilon$

$$< C_2 \epsilon^{R((a, \xi) \rightarrow (a', \xi'))}$$

→ Definition 8 [✓]

↪ explained

The stochastic potential of a state that is defined by  $(a, \xi)$  is the minimum total resistance of all

spanning trees that have state  $(a, \xi)$

as their 'root'  $\gamma(a, \xi) = \min_{T \in \mathcal{T}(a, \xi)} R(T)$

'stochastic potential'

set of all spanning trees rooted at  $(a, \xi)$

→ stochastic potential essentially indicates or

that quantifies how

'easy' it is to reach state

$(a, \xi)$  from all other states

in system.

Resistance of a path  $T$  is given by

$$R(T) = \sum_{(a, \xi) \rightarrow (a', \xi') \in T} R((a, \xi) \rightarrow (a', \xi'))$$

basically total resistance  $R(T)$  is the sum of resistance for all the edges in path  $T$ .

**Theorem 1**

→ relationship b/w SSE & Stochastic Potential  
↪ states with min (SP) or stochastic potential correspond to SSEs. These are the most resilient states requiring least effort to reach & hence dominate the sys.

→ Definition 9 [✓]

→ Corollary 1, 2, 3 [✓]



Pareto optimal Outcome; where no player can be made better off without making another player worse off. Such SSEs which result in 'POO' aren't necessarily NEs.

\* → SSEs that are both Nash equilibria & 'POO' or Pareto optimal are generally preferred

## Definition 10 & Algorithm 1

→ StGs generalize normal-form games by incorporating multiple stages, states & transitions b/w states.

→ goal is equilibrium selection for stochastic games

→ framework for the algo is modular, essentially builds up on  $K_E$

Key components:

→  $K_E$  governs Player's behavior; it determines how players adjust their actions & hidden var  $\xi$  at each iter based on rewards or value func.

→ Actor-Critic Structure:

→ Algo essentially has 2 main steps

Step 1

Actor

Updates player action  $a_h^{(t+1)}(s)$  & hidden vars  $\xi_h^{(t+1)}(s)$  using  $K_E$  (incorporating current state & rewards)

Step 2

Critic

Updates value function,  $Q_i, v_h^{(t+1)}(s_i, a)$  for each player using a Bellman-like iteration, integrating both immediate & expected future rewards

System is modular coz by using different  $K_E$ , we get different equilibrium selection results.

# Algo 1

- Value functions  $Q_{i,h}^{(0)}(s,a)$  initialized with immediate rewards  $r_{i,h}(s,a)$
- Terminal value  $V_{i,H+1}^{(t)}(s)$  is initialized to 0 for all players 'i', all stages 'h', all states 's' & all iterations 't'.
- Actions  $a_h^{(0)}(s)$  & hidden var  $\xi_h^{(0)}(s)$  are randomly initialized.
- [Initializations]

for all  $t, t \geq 0$  till total iterations;

for each stage  $h$  (in reverse order);

## Step 1 - Actor

- next action  $a_h^{(t+1)}(s)$  & hidden var  $\xi_h^{(t+1)}(s)$  sampled based on current values  $Q_{i,h}^{(t)}(s,a)$  using  $K_E$  (This step determines players policy at this stage & state)

## Step 2 - Critic

- value functions updated,  $V_{i,h}^{(t+1)}(s)$  &  $Q_{i,h}^{(t+1)}(s,a)$ ;

$$V_{i,h}^{(t+1)}(s) = \frac{1}{t+1} \sum_{\tau=1}^{t+1} Q_{i,h}^{(\tau)}(s, a_h^{(\tau)}(s))$$

→ avg of all past  $Q$  values at state 's'.

$$Q_{i,h}^{(t+1)}(s,a) = r_{i,h}(s,a) + \sum_{s'} P_h(s'/s, a) \underbrace{V_{i,h+1}^{(t+1)}(s')}_{\text{expected return from state } s}$$

↓  
state, action-value function