

Last Time

Simple Games, Nash Equilibrium

		defender	
		climb	descend
attacker	climb	(3, -5)	(-1, 0)
	descend	(0, -1)	(4, -4)

	$\pi(\text{climb})$	$\pi(\text{descend})$	U
attacker	0.375	0.625	1.5
defender	0.625	0.375	-2.5

This Time

Sequential Games

Zero Sum	$R^1(\bar{a}) = -R^2(\bar{a})$
Cooperative Game	$R^i(\bar{a}) = R(\bar{a})$
General Sum	

Markov Game
 (I, S, A, T, R, γ)
 \uparrow joint action

$$T(s' | s, \bar{a})$$

Stochastic / mixed policies might be needed to describe equilibria

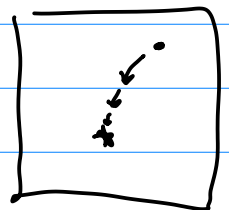
$$\pi^i(a | s) \quad \bar{\pi}(s)$$

Best Response

Fix policies of other agents ($\bar{\pi}^{-i}(s)$ is fixed)
 \hookrightarrow Results in MDP

Nash Equilibrium

Will Value Iteration Converge?



$$\begin{aligned} & \underset{\pi, U}{\text{minimize}} && \sum_{i \in I} \sum_s (U^i(s) - Q^i(s, \pi(s))) \\ & \text{subject to} && U^i(s) \geq Q^i(s, a^i, \bar{\pi}^{-i}(s)) \quad \forall i, s, a^i \\ & && \sum_{a^i} \pi^i(a^i | s) = 1 \quad \forall i, s \\ & && \pi^i(a^i | s) \geq 0 \quad \forall i, s, a^i \end{aligned}$$

where $Q^i(s, \bar{\pi}(s)) = R^i(s, \bar{\pi}(s)) + \gamma \sum_{s'} T(s' | s, \bar{\pi}(s)) U^i(s')$

Fictitious Play

$$\pi^i(a^i | s) \propto N^i(j, a^i, s)$$

Compute best response (solve MDP) periodically

Gradient Ascent

$$\frac{\partial U^{\pi_{t,i}}(s)}{\partial \pi_t^i(a^i | s)} = Q^{\pi_{t,i}}(s, a^i, a_{-i}^i)$$

Nash Q-learning

$\bar{Q}(s, \bar{a})$ joint action-value function

At every step

Compute Nash \bar{a}'

$$\bar{U}(s') = \sum_{\bar{a}'} \bar{Q}(s', \bar{a}') \prod_{j \in I} \pi^j(a^{j'})$$

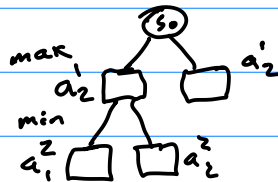
$$\bar{Q}(s, \bar{a}) \leftarrow \bar{Q}(s, \bar{a}) + \alpha (\bar{R}(s, \bar{a}) + \gamma \bar{U}(s') - \bar{Q}(s, \bar{a}))$$

$s \quad \bar{Q}(s, \bar{a})$

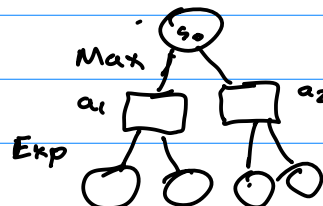
	a_1^z	a_2^z
a_1^i	\bar{Q}^i, \bar{Q}^z	\bar{Q}^i, \bar{Q}^z
a_2^i	\bar{Q}^i, \bar{Q}^z	\bar{Q}^i, \bar{Q}^z

Zero-Sum

Zero Sum Deterministic



MDP



State Uncertainty

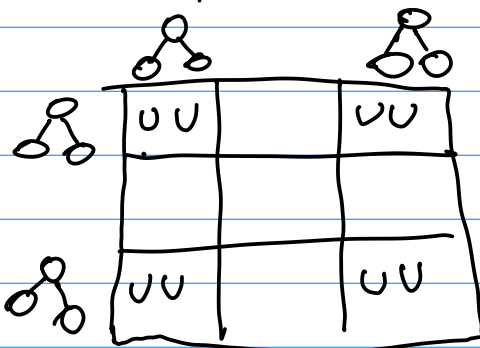
POMG

"POSG"

$(I, S, A, O, T, Z, R, \gamma)$

"Extensive-Form Game"

Belief updates "not possible"



Dynamic Programming

Similar to POMDP
value iteration

Differential Game

$$\dot{x} = f(t, x, u_1, \dots, u_N)$$

$$J_i = \int_0^T g_i(t, x, u_1, \dots, u_N) dt$$

Hamilton
Jacobi
Reachability

$$d \in D$$

Markov Game

$$x_{k+1} = G(x_k, \bar{a}, w)$$

