

**Decision Making**  
*in Robots and Autonomous Agents*

**Repeated, Stochastic and Bayesian Games**

Subramanian Ramamoorthy  
School of Informatics

26 February, 2013

# Repeated Game

- You can't learn if you only play a game once.
- Repeatedly playing a game raises new questions.
  - How many times? Is this common knowledge?

Finite Horizon

Infinite Horizon

- Trading off present and future reward?

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T r_t$$

Average Reward

$$\sum_{t=1}^{\infty} \gamma^t r_t$$

Discounted Reward

# Repeated Game - Strategies

- What can players do?
  - Strategies can depend on the history of play.

$$\sigma_i : \mathcal{H} \rightarrow PD(\mathcal{A}_i) \quad \text{where} \quad \mathcal{H} = \bigcup_{n=0}^{\infty} \mathcal{A}^n$$

- Markov strategies a.k.a. stationary strategies

$$\forall a^{1\dots n} \in \mathcal{A} \quad \sigma_i(a^1, \dots, a^n) = \sigma(a^n)$$

- $k$ -Markov strategies

$$\forall a_{1\dots n} \in \mathcal{A} \quad \sigma_i(a_1, \dots, a_n) = \sigma(a_{n-k}, \dots, a_n)$$

# Repeated Game - Examples

- Iterated Prisoner's Dilemma

$$R_1 = \begin{array}{c} \text{C} \\ \text{D} \end{array} \begin{array}{cc} \text{C} & \text{D} \\ \left( \begin{array}{cc} 3 & 0 \\ 4 & 1 \end{array} \right) \end{array} \quad R_2 = \begin{array}{c} \text{C} \\ \text{D} \end{array} \begin{array}{cc} \text{C} & \text{D} \\ \left( \begin{array}{cc} 3 & 4 \\ 0 & 1 \end{array} \right) \end{array}$$

- The single most examined repeated game!
- Repeated play can justify behavior that is not rational in the one-shot game.
- Tit-for-Tat (TFT)
  - \* Play opponent's last action (C on round 1).
  - \* A 1-Markov strategy.

# Well Known IPD Strategies

- AllC/D: always cooperate/defect
- Grim: cooperate until the other agent defects, then defect forever
- Tit-for-Tat (TFT): on 1<sup>st</sup> move, cooperate. On n<sup>th</sup> move, repeat the other agent's (n-1)<sup>th</sup> move
- Tit-for-Two-Tats (TFTT): like TFT, but only only retaliates if the other agent defects twice
- Tester: defect on round 1. If the other agent retaliates, play TFT. Otherwise, alternately C/D
- Pavlov: on 1st round, cooperate. Thereafter, win => use same action next; lose => switch

<i>AllC, Grim, TFT, or Pavlov</i>	<i>AllC, Grim, TFT, or Pavlov</i>	<i>TFT</i>	<i>Tester</i>	<i>Pavlov</i>	<i>AllD</i>
		C	D		
		D	C	C	D
C	C	C	C	D	D
C	C	C	C	C	D
C	C	C	C	D	D
C	C	C	C	C	D
C	C	C	C	D	D
⋮	⋮	⋮	⋮	C	D
				⋮	⋮

# Nash Equilibria – Repeated Game

- Obviously, Markov strategy equilibria exist.
- Consider iterated prisoner's dilemma and TFT.

$$R_1 = \begin{array}{c} \text{C} \\ \text{D} \end{array} \begin{array}{cc} \text{C} & \text{D} \\ \left( \begin{array}{cc} 3 & 0 \\ 4 & 1 \end{array} \right) \end{array} \quad R_2 = \begin{array}{c} \text{C} \\ \text{D} \end{array} \begin{array}{cc} \text{C} & \text{D} \\ \left( \begin{array}{cc} 3 & 4 \\ 0 & 1 \end{array} \right) \end{array}$$

- With average reward, what's a best response?
  - \* Always **D** has a value of 1.
  - \* **D** then **C** has a value of 2.5
  - \* Always **C** and TFT have a value of 3.
- Hence, both players following TFT is Nash.

# Nash Equilibria – Repeated Game

- The TFT equilibria is strictly preferred to all Markov strategy equilibria.
- The TFT strategy plays a dominated action.
- TFT uses a **threat** to enforce compliance.
- TFT is not a special case.

# Nash Equilibria – Repeated Game

**Folk Theorem.** For any repeated game with average reward, every *feasible* and *enforceable* vector of payoffs for the players can be achieved by some Nash equilibrium strategy. (Osborne & Rubinstein, 1994)

- A payoff vector is *feasible* if it is a linear combination of individual action payoffs.
- A payoff vector is *enforceable* if all players get at least their minimax value.



# Nash Equilibria – Repeated Game

**Folk Theorem.** For any repeated game with average reward, every *feasible* and *enforceable* vector of payoffs for the players can be achieved by some Nash equilibrium strategy. (Osborne & Rubinstein, 1994)

- Players' follow a deterministic sequence of play that achieves the payoff vector.
- Any deviation is punished.
- The threat keeps players from deviating as in TFT.

# Equilibria by ‘Learning’ – Universally Consistent

- A.k.a. Hannan consistent, regret minimizing.
- For a history  $h = a^1, a^2, \dots, a^n \in \mathcal{A}$ , define **regret** for player  $i$ ,

$$\text{Regret}_i(h) = \left( \max_{a_i \in \mathcal{A}_i} \sum_{t=1}^n R(\langle a_i, a_{-i}^t \rangle) \right) - \sum_{t=1}^n R_i(a^t)$$

i.e., the difference between the reward that could have been received by a stationary strategy and the actual reward received.

# Minimax by Regret Minimization

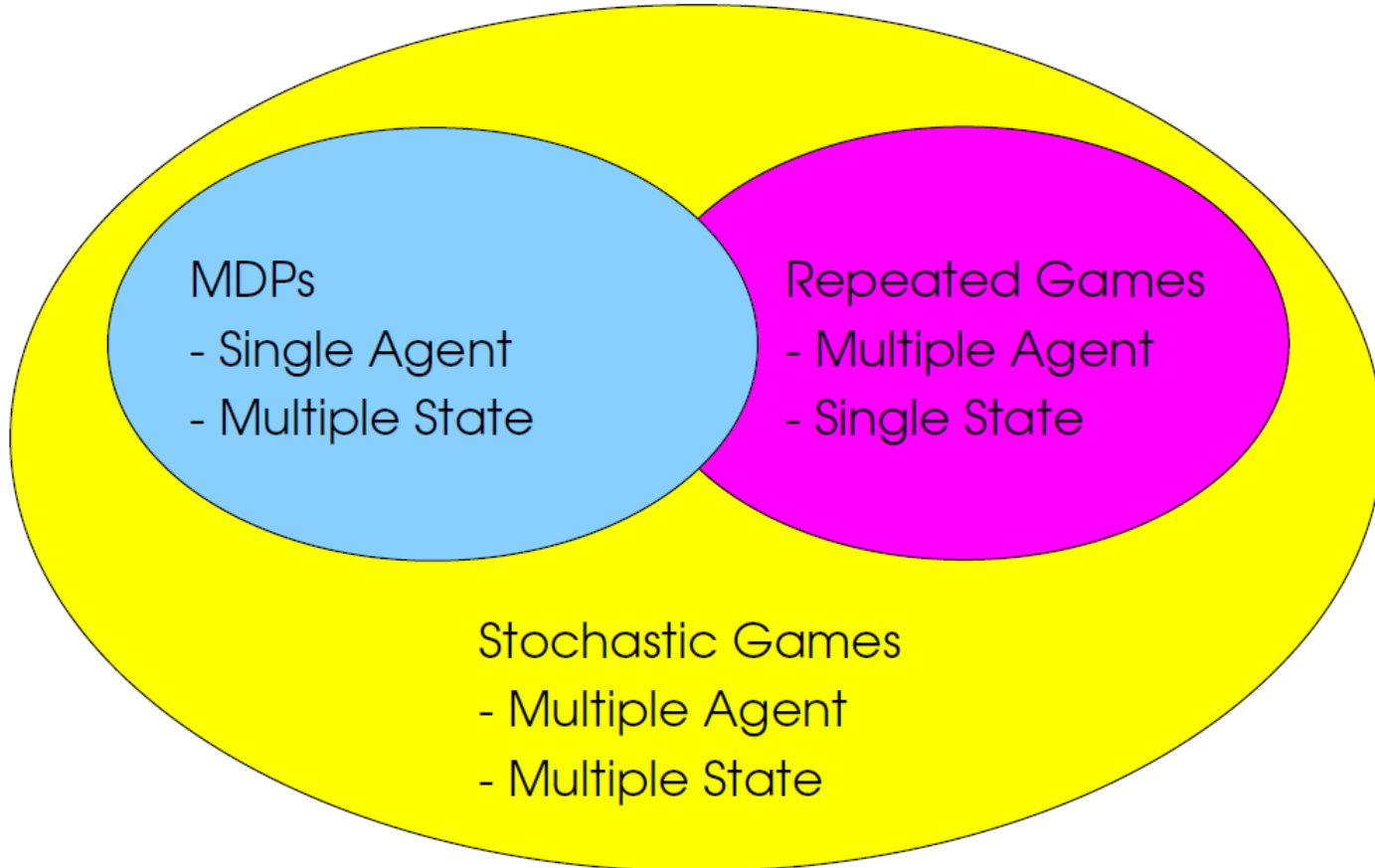
- A strategy  $\sigma_i$  is **universally consistent** if for any  $\epsilon > 0$  there exists a  $T$  such that for all  $\sigma_{-i}$  and  $t > T$ ,

$$\Pr \left[ \frac{\text{Regret}_i(a^1, \dots, a^t)}{t} > \epsilon \mid \langle \sigma_i, \sigma_{-i} \rangle < \epsilon \right]$$

i.e., with high probability the average regret is low for all strategies of the other players.

- If regret is zero, then must be getting at least the minimax value.

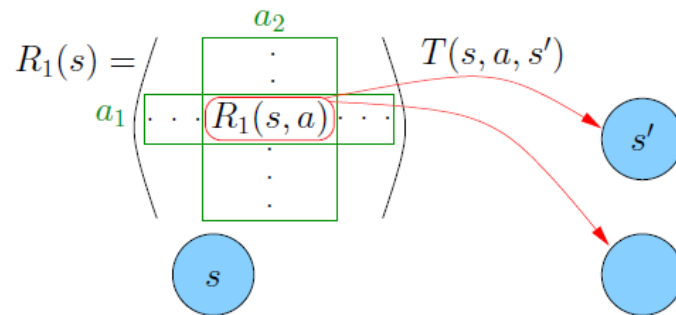
# Stochastic Games



# Stochastic Game - Setup

A stochastic game is a tuple  $(n, \mathcal{S}, \mathcal{A}_{1\dots n}, T, R_{1\dots n})$ ,

- $n$  is the number of agents,
- $\mathcal{S}$  is the set of states,
- $\mathcal{A}_i$  is the set of actions available to agent  $i$ ,
  - $\mathcal{A}$  is the joint action space  $\mathcal{A}_1 \times \dots \times \mathcal{A}_n$ ,
- $T$  is the transition function  $\mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ ,
- $R_i$  is the reward function for the  $i$ th agent  $\mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ .



# Stochastic Game - Policies

- What can players do?

- Policies depend on history and the current state.

$$\pi_i : \mathcal{H} \times \mathcal{S} \rightarrow PD(\mathcal{A}_i) \quad \text{where} \quad \mathcal{H} = \bigcup_{n=0}^{\infty} (\mathcal{S} \times \mathcal{A})^n$$

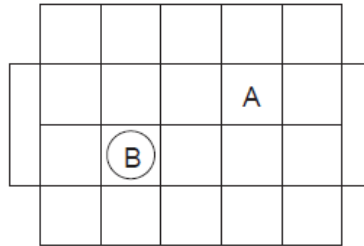
- Markov policies a.k.a. stationary policies

$$\forall h, h' \in \mathcal{H} \forall s \in \mathcal{S} \quad \pi_i(h, s) = \pi_i(h', s)$$

- Focus on learning Markov policies, but the learning itself is a non-Markovian policy.

# Stochastic Game - Example

(Littman, 1994)



- Players: Two.
- States: Player positions and ball possession (780).
- Actions: N, S, E, W, Hold (5).
- Transitions:
  - Simultaneous action selection, random execution.
  - Collision could change ball possession.
- Rewards: Ball enters a goal.

# Stochastic Game - Remarks

- If  $n = 1$ , it is an MDP.
- If  $|S| = 1$ , it is a repeated game.
- If the other players play a stationary policy, it is an MDP to the remaining player.

$$\hat{T}(s, a_i, s') = \sum_{a_{-i} \in \mathcal{A}_{-i}} \pi_{-i}(s, a) T(s, \langle a_i, a_{-i} \rangle, s')$$

- The interesting case, then, is when the other agents are not stationary, i.e., are learning.



# Nash Equilibria – Stochastic Game

- Consider Markov policies.
- A **best response set** is the set of all Markov policies that are optimal given the other players' policies.

$$\text{BR}_i(\pi_{-i}) = \left\{ \pi_i \mid \begin{array}{l} \forall \pi'_i \forall s \in \mathcal{S} \\ V_i^{\langle \pi_i, \pi_{-i} \rangle}(s) \geq V_i^{\langle \pi'_i, \pi_{-i} \rangle}(s) \end{array} \right\}$$

- A **Nash equilibrium** is a joint policy, where all players are playing best responses to each other.

$$\forall i \in \{1 \dots n\} \quad \pi_i \in \text{BR}_i(\pi_{-i})$$

# Nash Equilibria – Stochastic Game

- All discounted reward and zero-sum average reward stochastic games have at least one Nash equilibrium. (Shapley, 1953; Fink, 1964)

# Incomplete Information

- So far, we assumed that everything relevant about the game being played is common knowledge to all the players:
  - the number of players
  - the actions available to each
  - the payoff vector associated with each action vector
- True even for imperfect-information games
  - The actual moves aren't common knowledge, but the game is
- We'll now consider games of **incomplete (*not imperfect*) information**
  - Players are uncertain about the game being played

# Incomplete Information

- Consider the payoff matrix shown here
  - $\epsilon$  is a small positive constant; Agent 1 knows its value
- Agent 1 doesn't know the values of  $a, b, c, d$ 
  - Thus the matrix represents a **set of games**
  - Agent 1 doesn't know which of these games is the one being played
- Agent 1 seeks strategy that works despite lack of knowledge
- If Agent 1 thinks Agent 2 is malicious, then Agent 1 might want to play a maxmin, or “safety level,” strategy
  - minimum payoff of T is  $1-\epsilon$
  - minimum payoff of B is 1
- So agent 1's maxmin strategy is B

	$L$	$R$
$T$	$100, a$	$1 - \epsilon, b$
$B$	$2, c$	$1, d$

# Regret

- Suppose Agent 1 doesn't think Agent 2 is malicious
- Agent 1 might reason as follows:
  - If Agent 2 plays  $R$ , then 1's strategy changes 1's payoff by only a small amount
    - Payoff is 1 or  $1-\epsilon$ ;
    - Agent 1's difference is only  $\epsilon$
  - If Agent 2 plays  $L$ , then 1's strategy changes 1's payoff by a much bigger amount
    - Either 100 or 2, difference is 98
  - If Agent 1 chooses  $T$ , this will minimize 1's worst-case **regret**
    - Maximum difference between the payoff of the chosen action and the payoff of the other action

	$L$	$R$
$T$	100, $a$	$1 - \epsilon, b$
$B$	2, $c$	1, $d$

# Minimax Regret

- Suppose  $i$  plays action  $a_i$  and the other agents play action profile  $\mathbf{a}_{-i}$
- $i$ 's regret: amount  $i$  lost by playing  $a_i$  instead of  $i$ 's best response to  $\mathbf{a}_{-i}$

$$\text{regret}(a_i, \mathbf{a}_{-i}) = \left[ \max_{a'_i \in A_i} u_i(a'_i, \mathbf{a}_{-i}) \right] - u_i(a_i, \mathbf{a}_{-i})$$

- $i$  doesn't know what  $\mathbf{a}_{-i}$  will be, but can consider worst case:
  - maximum regret for  $a_i$ , maximized over every possible  $\mathbf{a}_{-i}$

$$\max_{\mathbf{a}_{-i} \in \mathbf{A}_{-i}} \text{regret}(a_i, \mathbf{a}_{-i}) = \max_{\mathbf{a}_{-i} \in \mathbf{A}_{-i}} \left( \left[ \max_{a'_i \in A_i} u_i(a'_i, \mathbf{a}_{-i}) \right] - u_i(a_i, \mathbf{a}_{-i}) \right)$$

# Minimax Regret

- **Minimax regret action:** an action with the smallest maximum regret

$$\operatorname{argmin}_{a_i \in A_i} \max_{\mathbf{a}_{-i} \in \mathbf{A}_{-i}} \operatorname{regret}(a_i) = \operatorname{argmin}_{a_i \in A_i} \max_{\mathbf{a}_{-i} \in \mathbf{A}_{-i}} \left( \left[ \max_{a'_i \in A_i} u_i(a'_i, \mathbf{a}_{-i}) \right] - u_i(a_i, \mathbf{a}_{-i}) \right)$$

- Can extend to a solution concept
  - All agents play minimax regret actions
  - This is one way to deal with the incompleteness, but often we can do more with the representation

# Bayesian Games

- In the previous example, we knew the set  $\mathbf{G}$  of all possible games, but didn't know anything about which game in  $\mathbf{G}$ 
  - Enough information to put a probability distribution over games
- A **Bayesian Game** is a class of games  $\mathbf{G}$  that satisfies two fundamental conditions
- **Condition 1:**
  - The games in  $\mathbf{G}$  have the same number of agents, and the same strategy space (set of possible strategies) for each agent. The only difference is in the payoffs of the strategies.
- This condition isn't very restrictive
  - Other types of uncertainty can be reduced to the above, by reformulating the problem



# An Example

- Suppose we don't know whether player 2 only has strategies L and R, or also an additional strategy C:

		<i>L</i>	<i>R</i>	
Game $G_1$	<i>U</i>	1, 1	1, 3	
	<i>D</i>	0, 5	1, 13	

		<i>L</i>	<i>C</i>	<i>R</i>	
Game $G_2$	<i>U</i>	1, 1	0, 2	1, 3	
	<i>D</i>	0, 5	2, 8	1, 13	

- If player 2 doesn't have strategy C, this is equivalent to having a strategy C that's strictly dominated by other strategies:
  - Nash equilibria for  $G_1'$  are the same as for  $G_1$

		<i>L</i>	<i>C</i>	<i>R</i>
Game $G_1'$	<i>U</i>	1, 1	0, -100	1, 3
	<i>D</i>	0, 5	2, -100	1, 13

- Problem is reduced to whether C's payoffs are those of  $G_1'$  or  $G_2$

# Bayesian Games

## ***Condition 2 (common prior):***

- The probability distribution over the games in  $\mathbf{G}$  is **common knowledge** (i.e., known to all the agents)
- So a Bayesian game defines
  - the uncertainties of agents about the game being played,
  - what each agent believes the other agents believe about the game being played
- The beliefs of the different agents are posterior probabilities
  - Combine the common prior distribution with individual “private signals” (what’s “revealed” to the individual players)
- The common-prior assumption rules out whole families of games
  - But it greatly simplifies the theory, so most work in game theory uses it

# The Bayesian Game Model

A Bayesian game consists of

- a set of games that differ only in their payoffs
- a common (known to all players) prior distribution over them
- for each agent, a partition structure (set of information sets) over the games

# Bayesian Game: Information Sets Defn.

A Bayesian game is a 4-tuple  $(N, G, P, I)$

- $N$  is a set of agents
- $G$  is a set of  $N$ -agent games
- For every agent  $i$ , every game in  $G$  has the same strategy space
- $P$  is a **common prior** over  $G$ 
  - **common**: common knowledge (known to all the agents)
  - **prior**: probability before learning any additional information
- $I = (I_1, \dots, I_N)$  is a tuple of partitions of  $G$ , one for each agent (information sets)

$G = \{\text{Matching Pennies (MP)}, \text{Prisoner's Dilemma (PD)}, \text{Coordination (Crd)}, \text{Battle of the Sexes (BoS)}\}$

	$I_{2,1}$	$I_{2,2}$												
$I_{1,1}$	MP ( $p = 0.3$ ) L R U <table border="1"><tr><td>2, 0</td><td>0, 2</td></tr><tr><td>0, 2</td><td>2, 0</td></tr></table> D <table border="1"><tr><td>0, 2</td><td>2, 0</td></tr></table>	2, 0	0, 2	0, 2	2, 0	0, 2	2, 0	PD ( $p = 0.1$ ) L R U <table border="1"><tr><td>2, 2</td><td>0, 3</td></tr><tr><td>3, 0</td><td>1, 1</td></tr></table> D <table border="1"><tr><td>3, 0</td><td>1, 1</td></tr></table>	2, 2	0, 3	3, 0	1, 1	3, 0	1, 1
2, 0	0, 2													
0, 2	2, 0													
0, 2	2, 0													
2, 2	0, 3													
3, 0	1, 1													
3, 0	1, 1													
$I_{1,2}$	Crd ( $p=0.2$ ) L R U <table border="1"><tr><td>2, 2</td><td>0, 0</td></tr><tr><td>0, 0</td><td>1, 1</td></tr></table> D <table border="1"><tr><td>0, 0</td><td>1, 1</td></tr></table>	2, 2	0, 0	0, 0	1, 1	0, 0	1, 1	BoS ( $p = 0.4$ ) L R U <table border="1"><tr><td>2, 1</td><td>0, 0</td></tr><tr><td>0, 0</td><td>1, 2</td></tr></table> D <table border="1"><tr><td>0, 0</td><td>1, 2</td></tr></table>	2, 1	0, 0	0, 0	1, 2	0, 0	1, 2
2, 2	0, 0													
0, 0	1, 1													
0, 0	1, 1													
2, 1	0, 0													
0, 0	1, 2													
0, 0	1, 2													

# Example

- Suppose the randomly chosen game is MP
- Agent 1's information set is  $I_{1,1}$ 
  - 1 knows it's MP or PD
  - 1 can infer **posterior probabilities** for each

$$\Pr[\text{MP}|I_{1,1}] = \frac{\Pr[\text{MP}]}{\Pr[\text{MP}] + \Pr[\text{PD}]} = \frac{0.3}{0.3 + 0.1} = \frac{3}{4}$$

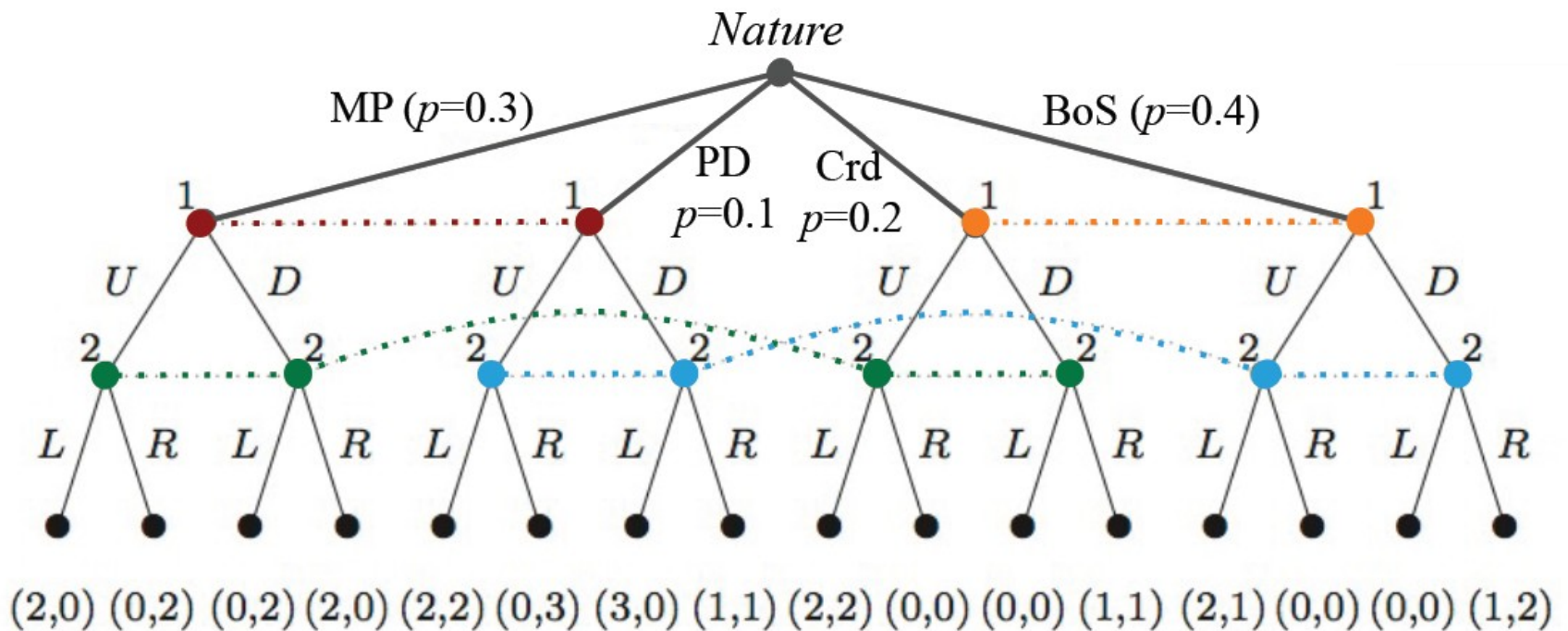
$$\Pr[\text{PD}|I_{1,1}] = \frac{\Pr[\text{PD}]}{\Pr[\text{MP}] + \Pr[\text{PD}]} = \frac{0.1}{0.3 + 0.1} = \frac{1}{4}$$

- Agent 2's information set is  $I_{2,1}$

$$\Pr[\text{MP}|I_{2,1}] = \frac{\Pr[\text{MP}]}{\Pr[\text{MP}] + \Pr[\text{CrD}]} = \frac{0.3}{0.3 + 0.2} = \frac{3}{5}$$

$$\Pr[\text{CrD}|I_{2,1}] = \frac{\Pr[\text{CrD}]}{\Pr[\text{MP}] + \Pr[\text{CrD}]} = \frac{0.2}{0.3 + 0.2} = \frac{2}{5}$$

# Another Interpretation: Extensive Form



# Epistemic Types

- We can assume the only thing players are uncertain about is the game's utility function
- Thus we can define uncertainty directly over a game's utility function

**Definition** : a **Bayesian game** is a tuple  $(N, A, \Theta, p, u)$  where:

$N$  is a set of agents;

$A = A_1 \times \dots \times A_n$ , where  $A_i$  is the set of actions available to player  $i$  ;

$\Theta = \Theta_1 \times \dots \times \Theta_n$ , where  $\Theta_i$  is the type space of player  $i$  ;

$p : \Theta \rightarrow [0, 1]$  is a common prior over types; and

$u = (u_1, \dots, u_n)$ , where  $u_i : A \times \Theta \rightarrow \mathfrak{R}$  is the utility function for player  $i$

- All this is common knowledge; each agent knows its own type

# Types

An agent's **type** consists of all the information it has that isn't common knowledge, e.g.,

- The agent's actual payoff function
- The agent's beliefs about other agents' payoffs,
- The agent's beliefs about *their* beliefs about his own payoff
- Any other higher-order beliefs



# Strategies

Similar to what we had in imperfect-information games:

- A pure strategy for player  $i$  maps each of  $i$ 's types to an action
  - what  $i$  would play if  $i$  had that type
- A mixed strategy  $s_i$  is a probability distribution over pure strategies

$$s_i(a_i|\theta_j) = Pr[i \text{ plays action } a_j \mid i \text{'s type is } \theta_j]$$

- Many kinds of expected utility: *ex post*, *ex interim*, and *ex ante*
  - Depend on what we know about the players' types

# Expected Utility

If we know every agent's type (i.e., the type profile  $\theta$ )

agent  $i$ 's *ex post* expected utility:

$$EU_i(\mathbf{s}, \theta) = \sum_{\mathbf{a}} \Pr[\mathbf{a} | \mathbf{s}, \theta] u_i(\mathbf{a}, \theta) = \sum_{\mathbf{a}} \left( \prod_{j \in N} s_j(a_j | \theta_j) \right) u_i(\mathbf{a}, \theta)$$

If we only know the common prior

agent  $i$ 's *ex ante*

expected utility:

$$EU_i(\mathbf{s}) = \sum_{\theta} \Pr[\theta] EU_i(\mathbf{s}, \theta) = \sum_{\theta_i} \Pr[\theta_i] EU_i(\mathbf{s}, \theta_i)$$

If we know the type  $\theta_i$  of one agent  $i$ , but not the other agents' types

$i$ 's *ex interim*

expected utility:

$$EU_i(\mathbf{s}, \theta_i) = \sum_{\theta_{-i}} \Pr[\theta_{-i} | \theta_i] EU_i(\mathbf{s}, (\theta_i, \theta_{-i}))$$

# Bayes-Nash Equilibrium

Given a strategy profile  $\mathbf{s}_{-i}$ , a **best response** for agent  $i$  is a strategy  $s_i$  such that

$$s_i \in \arg \max_{s'_i} (EU_i(s'_i, \mathbf{s}_{-i}))$$

Above, the set notation is because more than one strategy may produce the same expected utility

A **Bayes-Nash** equilibrium is a strategy profile  $\mathbf{s}$  such that for every  $s_i$  in  $\mathbf{s}$ ,  $s_i$  is a best response to  $\mathbf{s}_{-i}$

Just like the definition of a Nash equilibrium, except that we're using Bayesian-game strategies

# Computing Bayes-Nash Equilibria

- The idea is to construct a payoff matrix for the entire Bayesian game, and find equilibria on that matrix

	$\theta_{2,1}$	$\theta_{2,2}$												
$\theta_{1,1}$	MP ( $p = 0.3$ ) L R U <table border="1"><tr><td>2, 0</td><td>0, 2</td></tr><tr><td>0, 2</td><td>2, 0</td></tr></table> D <table border="1"><tr><td>0, 2</td><td>2, 0</td></tr></table>	2, 0	0, 2	0, 2	2, 0	0, 2	2, 0	PD ( $p = 0.1$ ) L R U <table border="1"><tr><td>2, 2</td><td>0, 3</td></tr><tr><td>3, 0</td><td>1, 1</td></tr></table> D <table border="1"><tr><td>3, 0</td><td>1, 1</td></tr></table>	2, 2	0, 3	3, 0	1, 1	3, 0	1, 1
2, 0	0, 2													
0, 2	2, 0													
0, 2	2, 0													
2, 2	0, 3													
3, 0	1, 1													
3, 0	1, 1													
$\theta_{1,2}$	Crd ( $p = 0.2$ ) L R U <table border="1"><tr><td>2, 2</td><td>0, 0</td></tr><tr><td>0, 0</td><td>1, 1</td></tr></table> D <table border="1"><tr><td>0, 0</td><td>1, 1</td></tr></table>	2, 2	0, 0	0, 0	1, 1	0, 0	1, 1	BoS ( $p = 0.4$ ) L R U <table border="1"><tr><td>2, 1</td><td>0, 0</td></tr><tr><td>0, 0</td><td>1, 2</td></tr></table> D <table border="1"><tr><td>0, 0</td><td>1, 2</td></tr></table>	2, 1	0, 0	0, 0	1, 2	0, 0	1, 2
2, 2	0, 0													
0, 0	1, 1													
0, 0	1, 1													
2, 1	0, 0													
0, 0	1, 2													
0, 0	1, 2													

Write each of the pure strategies as a list of actions, one for each type:

Agent 1's pure strategies:

- UU: U if type  $\theta_{1,1}$ , U if type  $\theta_{1,2}$
- UD: U if type  $\theta_{1,1}$ , D if type  $\theta_{1,2}$
- DU: D if type  $\theta_{1,1}$ , U if type  $\theta_{1,2}$
- DD: D if type  $\theta_{1,1}$ , D if type  $\theta_{1,2}$

Agent 2's pure strategies:

- LL: L if type  $\theta_{2,1}$ , L if type  $\theta_{2,2}$
- LR: L if type  $\theta_{2,1}$ , R if type  $\theta_{2,2}$
- RL: R if type  $\theta_{2,1}$ , L if type  $\theta_{2,2}$
- RR: R if type  $\theta_{2,1}$ , R if type  $\theta_{2,2}$

# Computing Bayes-Nash Equilibria

Compute *ex ante* expected utility for each pure-strategy profile:

$$\begin{aligned}
 EU_2(UU, LL) &= \sum_{\theta} \Pr[\theta] u_2(U, L, \theta) \\
 &= \Pr[\theta_{1,1}, \theta_{2,1}] u_2(U, L, \theta_{1,1}, \theta_{2,1}) \\
 &\quad + \Pr[\theta_{1,1}, \theta_{2,2}] u_2(U, L, \theta_{1,1}, \theta_{2,2}) \\
 &\quad + \Pr[\theta_{1,2}, \theta_{2,1}] u_2(U, L, \theta_{1,2}, \theta_{2,1}) \\
 &\quad + \Pr[\theta_{1,2}, \theta_{2,2}] u_2(U, L, \theta_{1,2}, \theta_{2,2}) \\
 &= 0.3(0) + 0.1(2) + 0.2(2) + 0.4(1) \\
 &= 1
 \end{aligned}$$

	$\theta_{2,1}$	$\theta_{2,2}$	
	MP ( $p = 0.3$ )	PD ( $p = 0.1$ )	
	L R	L R	
$\theta_{1,1}$	U	2, 0	0, 2
	D	0, 2	2, 0
	Crd ( $p = 0.2$ )	BoS ( $p = 0.4$ )	
	L R	L R	
$\theta_{1,2}$	U	2, 2	0, 0
	D	0, 0	1, 1

# Computing Bayes-Nash Equilibria

- Put all of the *ex ante* expected utilities into a payoff matrix

e.g.,  $EU_2(UU,LL) = 1$

- Now we can compute best responses and Nash equilibria

		$\theta_{2,1}$		$\theta_{2,2}$	
		MP ( $p = 0.3$ )		PD ( $p = 0.1$ )	
		L	R	L	R
$\theta_{1,1}$	U	2, 0	0, 2	2, 2	0, 3
	D	0, 2	2, 0	3, 0	1, 1
		Crd ( $p=0.2$ )		BoS ( $p = 0.4$ )	
		L	R	L	R
$\theta_{1,2}$	U	2, 2	0, 0	2, 1	0, 0
	D	0, 0	1, 1	0, 0	1, 2

	<i>LL</i>	<i>LR</i>	<i>RL</i>	<i>RR</i>
<i>UU</i>	2, 1	1, 0.7	1, 1.2	0, 0.9
<i>UD</i>	0.8, 0.2	1, 1.1	0.4, 1	0.6, 1.9
<i>DU</i>	1.5, 1.4	0.5, 1.1	1.7, 0.4	0.7, 0.1
<i>DD</i>	0.3, 0.6	0.5, 1.5	1.1, 0.2	1.3, 1.1

# Computing Bayes Nash Equilibria

- Suppose we learn agent 1's type is  $\theta_{1,1}$   
 Recompute the payoff matrix using the posterior probabilities  
 $\Pr[\text{MP}|\theta_{1,1}] = \frac{3}{4}$ ,  $\Pr[\text{PD}|\theta_{1,1}] = \frac{1}{4}$
- $u_2(UU, LL|\theta_{1,1}) = \frac{3}{4}(0) + \frac{1}{4}(2) = 0.5$
- *Ex interim* payoff matrix when agent 1's type is  $\theta_{1,1}$
- Can't use this to compute equilibria, because  $\theta_{1,1}$  isn't common knowledge

		$\theta_{2,1}$		$\theta_{2,2}$	
		MP ( $p = 0.3$ )		PD ( $p = 0.1$ )	
		L	R	L	R
$\theta_{1,1}$	U	2, 0	0, 2	2, 2	0, 3
	D	0, 2	2, 0	3, 0	1, 1

	<i>LL</i>	<i>LR</i>	<i>RL</i>	<i>RR</i>
<i>UU</i>	2, 0.5	1.5, 0.75	0.5, 2	0, 2.25
<i>UD</i>	2, 0.5	1.5, 0.75	0.5, 2	0, 2.25
<i>DU</i>	0.75, 1.5	0.25, 1.75	2.25, 0	1.75, 0.25
<i>DD</i>	0.75, 1.5	0.25, 1.75	2.25, 0	1.75, 0.25

and so on...

# Acknowledgements

Slides are adapted from:

- Tutorial at IJCAI 2003 by Prof Peter Stone, University of Texas
- Game Theory lectures by Prof. Dana Nau, University of Maryland