

SUMMARY OF AAMAS

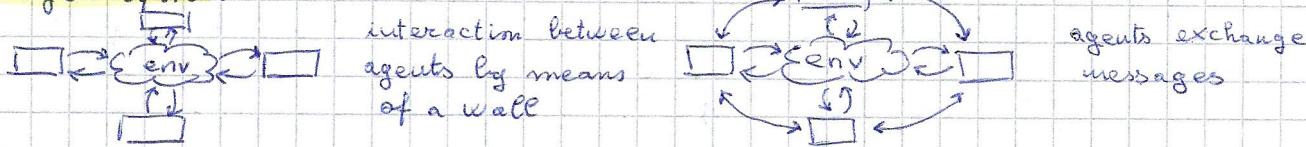
Autonomous agent: how can it decide his actions? \rightarrow to maximize an utility function

- PREPROGRAM: decided by designer, environment very controlled

- ONLINE DECISION: equip the agent with the ability of decide what to do

- LEARNING AGENTS: equip with the ability of learn, improve performances over time.

Multagent system:



Interaction:

- cooperative: to reach the global goal of the system, interaction to do sth together
- competitive: agents try to maximize proper utility (no general goal)

Regulate behavior

- pre-programmed decisions: efficient but doesn't cope with dynamic environment or errors
- autonomous behavior: you tell the goal and leave freedom but you have to find a way to regulate the behavior of the single component

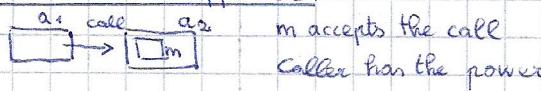
Human Organization

- entities: heterogeneous, cognitive (humans) or autonomous (agents)
- relationships: create an order of behaviors (es: A boss of B, negotiations)
- norms: constraints and regulation that cannot be violated (sth. you have a cost)
- environment: the more complex it is, the more complex the system will be.

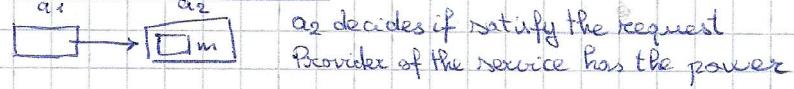
Agent properties:

- autonomy: operate without human supervision
- social ability: interaction with other agents, humans and environment
- reactivity: respond in a fast way to changes in the environment
- proactivity: goal directed behavior Composition of reactive and proactive behavior

Client-server approach



Agent-based approach



- reasoning agents (AI): agent able to perceive sth, decide and execute (planning, STRIPS, PDDL, A*)

- reactive agents: global behavior obtained by composing simple behaviors (subsumption architecture), when several behaviors are applicable define an order of priority.

- BDI agents: believes (what you know), desires (goals), intention (plan to execute)

- Layered agents: usually reactive in low, reasoning going higher.

Three levels: reactive-planning (pro-active)- communication (social)

MARKOV DECISION PROCESSES (MDPs)

S: states that the agent can distinguish

A: action " " is able to perform

(transition function)

P: $S \times A \rightarrow \Delta S$ set of all prob distribution over all elements of S. $P(s'|s,a)$ $S \times A \times S \rightarrow [0,1]$

r: $S \times A \rightarrow R$ reward function (gain)

Policy $d: S \rightarrow A$ decision, not care the effect.

If I knew the next state: $s \xrightarrow{d(s)=a} s' \xrightarrow{d(s')=a'} s'' \dots \max r(s,a) + r(s',a') + r(s'',a'') + \dots$

problems: unknown next state AND absence of terminal state (no reward with any policy)

Discount factor λ ($0 < \lambda < 1$): $r(s,a) + \lambda r(s',a') + \lambda^2 r(s'',a'') + \dots$

Value function $v^d(s)$: expected reward starting from s and following policy d

Bellman equation

Optimal policy d^* (you don't know which one) has $v^*(s) = \max_{a \in A} [r(s,a) + \lambda \sum_{s' \in S} p(s'|s,a) v^*(s')]$

$d^*(s) = \arg \max_{a \in A} [r(s,a) + \lambda \sum_{s' \in S} p(s'|s,a) v^*(s')]$ \rightarrow system of non linear equations for $v^*(s)$

Value iteration $v^* \leftarrow \text{random initialization}$ // always converge

algorithm: repeat $v \leftarrow v^*$, for each $s \in S$

$$v^*(s) \leftarrow \max_{a \in A} [r(s,a) + \lambda \sum_{s' \in S} p(s'|s,a) v(s')]$$

until $\max |v(s) - v^*(s)| < \epsilon$, return v^*

NON-COOPERATIVE GAMES

Mechanism (rules): NORMAL-FORM-GAMES \rightarrow players play simultaneously

N : set of players

$A = \{A_1, A_2, \dots, A_n\}$ A_i : set of actions of player i

$U_i \in \{U_1, U_2, \dots, U_n\}$ $U_i: A_1 \times A_2 \times \dots \times A_n \rightarrow \mathbb{R}$ utility function

	R	P	S
Ex: Rock-Paper-Scissors	R: 0,0	-1,1	1,-1
$N = \{1, 2\}$	P: 1,-1	0,0	-1,1
$A_i = \{R, P, S\}$	S: -1,1	1,-1	0,0

Strategies:

Behavior of players. In principle, a f that returns ~~one action~~ but we cannot bluff: we need distribution

$$S_1 = \left\{ \begin{array}{l} \frac{1}{3} R \\ \frac{1}{3} P \\ \frac{1}{3} S \end{array} \right\}$$

$$S_2 = \left\{ \begin{array}{l} \frac{1}{6} R \\ \frac{1}{2} P \\ \frac{1}{3} S \end{array} \right\} \rightarrow \text{Lottery probability distribution over outcomes (expected utility)}$$

$$E[U_1(S_1, S_2)] = 0 + \frac{1}{6} - \frac{1}{6} + 0 + \frac{1}{6} = 0$$

\leq : strictly dom.

\leq : weakly dom

Strategy is not good in all the possible situations. Properties of the strategies:

\rightarrow DOMINANCE: given a player i and a $a \in A_i$, a is dominated IF $\exists a' \in A_i: U_i(a_i, a_{-i}) \leq U_i(a'_i, a_{-i})$ $\forall a_j$ (a_i always worse, for each possible actions made by the opponents).

IF after discarding dominated a it remains one action per player, dominant strategy equilibrium. T 10/20

Assumption: I know utility of the opponent \rightarrow Iterated dominance (I can discard further actions C 21/34 assuming what a player will play). Not always dominance can be applied.

\rightarrow NASH EQUILIBRIUM: (S_1, S_2) is a NE IF $\forall i \in N, E[U_i(S_1, S_2)] \geq E[U_i(S_i^*, S_{-i})]$ VS:

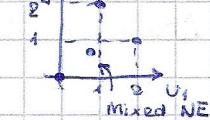
Check that players don't deviate. In mixed strategy players are randomized, otherwise it could be No equilibrium \rightarrow the opponent understand if you always play the same.

Nash proved that in every finite game there is at least a Nash Equilibrium (guaranteed in mixed strategies) and the # of equilibria is odd.

Player is rational: randomization introduces a constraint

$$\begin{cases} S_1(B) > 0 \\ S_2(S) > 0 \end{cases} \Rightarrow E[U_1(B, S_2)] = E[U_1(S, S_2)] \Rightarrow S_2(B) = 1 \cdot (1 - S_2(B)) \Rightarrow S_2(B) = \frac{1}{3}$$

Both or Starvinsky game



Nash Equilibrium (and also use dominance) can be inefficient \rightarrow rational players are selfish.

\rightarrow SUBGAME PERFECT EQUILIBRIUM: it is a Nash Equilibrium of every subgame of the original game.

Every finite extensive game has a subgame perfect eq.

\hookrightarrow it allows to represent the sequence of players' moves and the payoffs, we have:

- perfect info game: every player observes all the actions of the other players;

- imperfect info game: player cannot observe action of another player, we link possible resulting nodes into an information set

Propagate backward in the tree the decisions of the players, substitute an incomplete tree with his Nash equilibrium to find the subgame perfect equilibrium.

NEGOTIATION: agents with diff interests try to reach an agreement without supervisors (peers)

Agents

Issues: objects under negotiation (e.g. car)

Outcomes: results of negotiation (e.g. price)

O: set of outcomes (e.g. all the possible #'s)

$U^i: O \rightarrow \mathbb{R}$: utility of agent i

protocol \rightarrow rules defining the negotiation (e.g. low the price if there are more convenient offers)

strategy \rightarrow individual, how agents act in the protocol (e.g. how to low it)

- cooperative (axiomatic) approach: fix some axioms to predict the outcome;

- non-cooperative (algorithm) appr: agents are free, run an algorithm to discover the outcome.

\rightarrow Axiomatic approach: agents perfectly rational (always prefer an outcome that gives higher reward).

$O = A \cup \{D\}$ A: set of agreements D: disagreement $i \in \{A, D\} \rightarrow$ the two agents

bargaining set $S = \{(U^a(x), U^b(x)) | x \in A\}$

$d = (d^a, d^b) = (U^a(D), U^b(D))$

one element from the domain

bargaining problem (S, d) , bargaining solution $f: (S, d) \mapsto (x^a, x^b) \in S$

S closed and convex: $(x^a, x^b) \in S, (y^a, y^b) \in S \rightarrow \theta x^a + (1-\theta)y^a \in S \forall \theta \in [0, 1]$

$S \cap \{(x^a, x^b) | x^a \geq d^a, x^b \geq d^b\} \neq \emptyset$

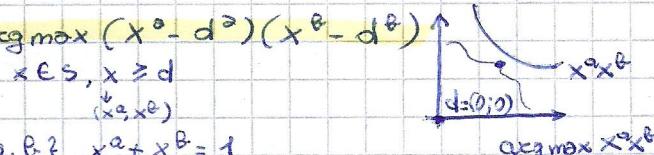
Properties of f :

- individual rationality: $f(S, d) \geq d$

- symmetry: if I swap A and B, I obtain the same point (just swapped)

- pareto optimality: impossible to improve it without making at least one agent

- no $y \in S$ s.t. $y \geq f(s, d)$ $\forall i$
- invariance: $J^a = \alpha x^a + \beta$, solution invariant to transformations (just stretched) $f(s, d) \in S'$
 - independence of irrelevant alternatives: $\forall S' \subseteq S$ (S' convex and containing the solution) $f(S', d) = f(S, d)$
- Given this, Nash bargaining solution: $f(S, d) \in \arg\max_{x \in S, x \geq d} (x^a - d^a)(x^b - d^b)$



\rightarrow Rubinstein's alternating offer protocol agents $\{a, b\}$ $x^a + x^b = 1$
 $O = \{(x^a, x^b) \mid 0 \leq x^a, x^b \leq 1, x^a + x^b = 1\}$ $t = 1, 2, \dots$ at each time agent i does an offer

$$U^a(x^a, t) = x^a \delta_a^{t-1} \quad 0 \leq \delta_a \leq 1 \rightarrow \text{discount factor}$$

$\delta_a = 1$: agent A patient

$\delta_a \approx 0$: impatient, after some point utility is 0.

$$U^b(x^b, t) = x^b \delta_b^{t-1}$$

We can have deadline or not, without agreement, utility = 0.

$$\text{If negotiation can last forever: } x^a = \frac{1 - \delta_b}{1 - \delta_a \delta_b}; \quad x^b = \frac{\delta_b - \delta_a \delta_b}{1 - \delta_a \delta_b} \quad \begin{cases} \delta_a = 1 \rightarrow x^a = 1, x^b = 0 \\ \delta_b \approx 0 \end{cases}$$

$n=1$ Ag A propose $(1; 0)$, Ag B accepts (or no agreement)

$n=2$ Ag A propose δ that Ag B can accept (\geq to the MAX that ag B can obtain at $t=2 \rightarrow \delta$)
 $\hookrightarrow (1 - \delta, \delta)$. Following this strategy the first offer is accepted.

For $t < N$ ag i OFFER $(1 - \delta x^j(t+1), \delta x^j(t+1))$ at $t=N$ ag i OFFER $(1; 0)$ if it's his turn
otherwise accepts the OFFER from j .

The protocol holds when $V^i = x^i \delta_i^{t-1}$ and the resource is divided between i and j.

\rightarrow Monotonic concession protocol: the two agents make an offer at the same time

$$(1) x^{(a)} \leftarrow \arg\max_{x \in O} U^a(x)$$

(2) propose $x^{(a)}$ to agent B

(3) receive proposal $x^{(b)}$ from ag B

(4) if $U^a(x^{(b)}) \geq U^a(x^{(a)})$ then accept $x^{(b)}$ and return

else $x^{(a)} \leftarrow x \in O$ st. $U^{a+}(x) \geq U^{b+}(x^{(a)})$ (and $U^a(x) \geq 0$)

(5) go to (2)

$$U^b(x) = U^b(x^{(a)}) + \epsilon$$

1) easy to verify if the other agent follows the protocol (concesses not)

2) the protocol converges to a solution but it can be very slow

not always you know U of the other (to make a concession).

Variant (concede in proportion to how much they are conceding): Zeuthen strategy

(in the else branch) we compute the willingness to risk a breakdown:

$$\text{risk}_a = \frac{U^a(x^{(a)}) - U^a(x^{(b)})}{U^a(x^{(a)})}$$

if $\text{risk}_a < \text{risk}_b$

then $x^{(a)} \leftarrow x \in O$ st. $\text{risk}_a(x) > \text{risk}_b(x)$ and goto(2)

$$\text{risk}_b = \frac{U^b(x^{(b)}) - U^b(x^{(a)})}{U^b(x^{(b)})}$$

else goto(3)

It converges quickly. The solution is a Nash bargaining solution.

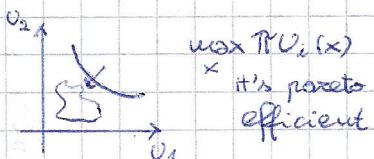
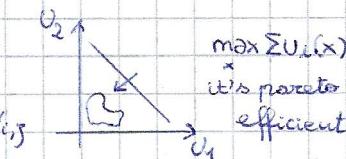
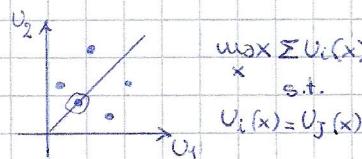
Solution concepts:

• PARETO EFFICIENCY

• EQUALITARIAN SOL

• UTILITARIAN SOL

• NASH BARGAINING SOL



COMPUTATIONAL SOCIAL CHOICE

agents $N = \{1, 2, \dots, n\} \rightarrow$ each one is a voter

alternatives $U \quad |U| \geq 2$

pref relation $\succ_i: a \succ_i b \rightarrow a$ is preferred than b by agent i

pref profile $R = [\succ_1, \succ_2, \dots, \succ_n]$

$R(U)$: set of all possible pref relation $\# = |U|^{|U|}$

$R(U)^n$: " " vectors of pref profile

-transitive: $a \succ_i b, b \succ_i c \rightarrow a \succ_i c, a, b, c \in U$

-complete: $a \succ_i b$ or $b \succ_i a$ or both

Vote means aggregate preferences and come up with an order

Aggregation function: social welfare function $f: R(U)^n \rightarrow R(U) \quad f([\succ_1, \succ_2, \dots, \succ_n]) = \succ$

1) Pareto optimality: if $a >_i b \forall i \in N$ then $a > b$ (strict preference)

2) independence of irrelevant alternatives (IIA): changing positions in pref relation, but NOT the order for a and b → obtain a diff final order BUT the same relation b/w a, b

$R[\{a, b\}]$: restriction of pref profile applied to a and b , $f(R) = \succ$

if $R[\{a, b\}] = R'[\{a, b\}]$ then $\succ[\{a, b\}] = \succ'[\{a, b\}]$

3) non-dictatorship: $\exists i \in N$ s.t. $a \succ_i b \rightarrow a \succ b$ for any R and $a, b \in U$.

Arrow's theorem: there is no f s.t. it satisfies these three properties (IIA usually drops).

Social choice function $f: R(U)^n \times \mathcal{P}(U) \rightarrow \mathcal{P}(U) \quad \rightarrow$ set of all possible subsets of feasible alternatives.

$A \in \mathcal{P}(U) \quad f([\succ_1, \succ_2, \dots, \succ_n], A) \rightarrow A' \quad A' \subseteq A$. How to implement it?

→ Voting rule $f: R(U)^n \rightarrow \mathcal{P}(U)$ f is resolute if $|f(R)| = 1$

→ Scoring systems: family of voting rules, score every alternative.

$\langle s_1, s_2, \dots, s_m \rangle \quad m = |U| \quad s_1 \geq s_2 \geq \dots \geq s_m, s_1 > s_m$

- plurality voting: $\langle 1, 0, \dots, 0 \rangle$

- Borda count: $\langle n-1, n-2, \dots, 1, 0 \rangle$

Alternative preferred by most of the agents

Condorcet winner: alternative that wins against each other alternative with pairwise majority

→ Condorcet extension: voting system that selects the condorcet winner (if it exists)

- Copeland's rule: an alt gets one point every time it wins ($\frac{1}{2}$ in case of tie) → cond extension

→ Multistep voting systems:

- Plurality with runoff: take the two most voted alternatives from a plurality voting and run a pairwise majority. Choosing how to vote has a large influence on the outcome

- Single Transferable Vote (STV): run $n-1$ plurality voting, excluding the last alt every time

- Sequential Majority Voting: sequence of pairwise majority voting following an Agenda (order) → pairwise maj between the first two, the winner goes against the third alt. $\hookrightarrow \# = \frac{n(n-1)}{2}$

Every agent has the model of preferences of the others: strategic manipulation of voting to avoid dissatisfaction

Gibbard-Satterthwaite Th: every possible reasonable voting system is prone to manipulation: it's

possible to represent a not-true order of preferences in order to have a winner that is better for the voting agent. (easy to estimate pref relation of the others).

→ voting systems for which manipulation is computationally high (cryptography)

→ Probabilistic voting systems: introduce random elements to force the agents to tell the truth.

Th: no voting system that ① avoid manipulation, ② is pareto efficient and ③ not probabilistic dictator

AUCTIONS: different agents make a bid, an auctioneer will take a decision according to a mechanism design

single-good auction → one single item is auctioned

multiple-good auction → more tricky, you can submit bids for multiple goods (combinatorial auctions)

Scarce resource to allocate on only one of the agents:

- English: auctioneer fixes an initial price, agents increase it until there won't be an higher offer, then the product will be sold. Problems: winner determination problem and payment.

- Japanese: auctioneer sets an initial price and increase it, meanwhile agents can leave the room. The remaining agent will be the winner and it will pay the actual value.

- Dutch: dual of Japanese → auctioneer starts from a very high price and decreases it until an agent accepts the offer. Control of time: you know the maxtime. Only one signal to stop.

- Sealed-bid: every agent makes a bid in a closed envelope. It's a single decision (diff from previous)

• first price: the auctioneer will choose the agent who submitted the higher bid

• second price: winner is the larger bid, the agent will pay the 2nd bid. (Vickrey)

Assumption: agents have independent private values (evaluation interp. drawn from a distribution, false in real world). Dominant strategy (in second price situation): bid exactly v_i .

You reveal your v_i , this problem is not present in first price.

In English and Japanese auctions the dominant strategy is the same, not true for Dutch or first price auctions \rightarrow no dominant strategy but equilibrium.

Th: n agents risk-neutral (evaluation without consider risk), independent v_i for each agent, then the equilibrium strategy is: every agent bids $\frac{n-1}{n} v_i$ (you use it if you know the others use it).

If n agents risk-neutral with independent v_i , interval of values $[v, \bar{v}]$ and:

① assign the good to the agent who offers more ② never assign to the agent with $v_i = \bar{v}$ then, every mechanism brings to the auctioneer the same expected revenue.

COMBINATORIAL AUCTIONS

agents $N = \{1, 2, \dots, n\}$

goods $G = \{1, 2, \dots, m\}$

possible subsets of $G: 2^m$ impossible to represent.

valuation function $v_i: 2^G \rightarrow \mathbb{R} \quad i \in N \quad \hat{v}_i: \text{bid}$

- Additive: $v_i(S \cup S') = v_i(S) + v_i(S')$ assuming $S \cap S' = \emptyset$. Specify v_i of single elements
- Complementary: $v_i(S \cup S') > v_i(S) + v_i(S')$ additional reward (e.g.: shoes)
- Substitutive: $v_i(S \cup S') < v_i(S) + v_i(S')$ dual case of the previous

Auctioneer allocates goods to the agents in order to maximize the revenue.

Allocation x : specify which subset of G goes to every agent (partition of G)

X : set of possible allocations ($x \in X$).

Goal of auctioneer: $\max \sum_{x \in X} \sum_{i \in N} \hat{v}_i(s) x_{s,i} \quad x_{s,i}: \text{Binary, if set } s \text{ is allocated to } i$

Impositions: $\sum_{\substack{s \in G \\ \text{size } s \leq 1}} x_{s,i} \leq 1 \quad \forall j \in N; \sum_{s \in G} x_{s,i} \leq 1 \quad \forall i \in N; x_{s,i} \in \{0, 1\} \quad \forall s \in G \quad \forall i \in N$
(technical const.)

Source of complexity: $n \times m$ variables \rightarrow for large instances no optimal solution.

For every set keep track only of $\bar{v}(s)$: max offer received for s . We add $\bar{v}(s) = 0$ for single items if they are missing. For non-single items apply additivity in case of missing.

Branch on items: build a tree of possible allocations, starting from the empty one. At each step, I add the bids containing the i -th good but not goods already in the path.

You can calculate the utility of each path and take the higher revenue. (keep track of the agent that does each offer \bar{v}). To speed up:

$x^* \leftarrow \{\}$

branch-on-items(j, x)

for sets S s.t. $j \in S$ and $S \cap (U_S \setminus x) = \{\}$

$g^* \leftarrow \emptyset$

$x' \leftarrow x + S$

branch-on-items($1, \{\}$)

if sets S in x contain all goods in G

if $\sum_{s \in x} \bar{v}(s) > g^*$

then if $\sum_{s \in x} \bar{v}(s) > g^*$

$\sum_{s \in x} \bar{v}(s) + h(x') > g^*$

then $g^* \leftarrow \sum_{s \in x} \bar{v}(s)$

then branch-on-items(j', x')

$x^* \leftarrow x$

return

Tree built in a depth-first fashion.

$h(x')$: estimate of the value I can get from allocating objects not yet in x' (upper bound)

$\hookrightarrow = \sum_{\substack{\text{several} \\ \text{not in } x'}} \max_{S \in G} \bar{v}(s) \quad \text{optimistic estimate}$

... again on auctions

Utility: $U(x) = RP_i - x$ where RP =reservation price, x : what you pay ($x \leq RP_i$)

Allocation function: f that selects the winner

On first price $v_i = 0$, in second price the winner has $v_i > 0$.

(1st p) \hookrightarrow winner decreases the price until the second bid to obtain a NE, all the players should know params on the others \rightarrow NE is not used to prescribe a strategy.

(2nd p) \hookrightarrow weakly dominant strategy in telling the truth (best you can do independently from the others, you don't need other info)

COALITION FORMATION: cooperative games, each agent has benefit in cooperating and can be bound in an agreement (constraint)

Transferable Utility Games (class of cooperative games): utility/payoff is given to the coalition, then the team decides how to divide the amount between the members. Utility given to a team is independent from the utility given to other teams (not in football).

$A = \{a_1, \dots, a_n\}$ set of agents

$V: 2^A \rightarrow \mathbb{R}$ characteristic function. The returned value represents how well a subset of agents is performing a task. $V(\emptyset) = 0 \quad \forall c \subseteq A$ A : largest coalition, grand coalition
 V is the value of coalition, not how the agents will divide it.

Outcome (CS, \vec{x}) \Rightarrow coalition structure \rightarrow partition of the set of agents (set of coalitions)

$$CS = \{C_1, \dots, C_m\} \text{ s.t. } \bigcup_i C_i = A \quad C_i \cap C_j = \emptyset \quad i \neq j$$

\vec{x} : payoff vector $\rightarrow \vec{x} = (x_1, \dots, x_n)$, $x_i \geq 0 \quad \forall i$; x_i is the payoff for ag; we assume that $\sum_i x_i = V(C)$

On the family of transferable utility games:

- super additive game: $V(C_1 \cup C_2) \geq V(C_1) + V(C_2) \quad \forall C_1, C_2 \subseteq A$ and disjointed Best CS: gran coalition
- convex game: $V(C_1 \cup C_2) + V(C_1 \cap C_2) \geq V(C_1) + V(C_2) \quad \forall C_1, C_2 \subseteq A$
they are super additive (taken into account for disjointed subsets), more special
- simple game: $V(C) \in \{0, 1\} \quad \forall C \subseteq A$

How to find an interesting outcome? Two steps: finding a good coalition structure and a good payoff vector. Payoff vector is NOT stable if a subset of a coalition has a sum of payoffs that is lower than the payoff of their coalition (unhappiness)

CORE: set of all stable outcomes. $Core(G) = \{(CS, \vec{x}) \mid \sum_{i \in C} x_i \geq V(C) \quad \forall C \subseteq A\}$

Problems: several outcomes in the core seem to be unfair

- $A = \{1, 2, 3\} \quad V(C) = \begin{cases} 1 & \text{if } |C| > 1 \\ 0 & \text{otherwise} \end{cases}$ No possibility to distribute 1 keeping stability. \rightarrow empty core (not in convex games)

Shapley values: idea is to produce a fair outcome, based on marginal contribution (in super additive) permutation $\{a_1, a_2, \dots, a_n\} = A$. $n!$ possible permutations $\Pi \in \Pi_A$ - set of all possible permutations $C_{\Pi}(a_i) \rightarrow$ set of agents that appear before a_i in permutation Π

$\Delta_{\Pi}^G(a_i) = V(C_{\Pi}(a_i) \cup \{a_i\}) - V(C_{\Pi}(a_i)) \rightarrow$ marginal contribution

$$\varphi_i(G) = \frac{1}{n!} \sum_{\Pi \in \Pi_A} \Delta_{\Pi}^G(a_i) \rightarrow$$
 Shapley value. Payoff vector composed by Shapley values.
 $CS = A$ (super additive game)

Properties:

1) Efficiency: $\varphi_1(G) + \varphi_2(G) + \dots + \varphi_n(G) = V(A)$, $\vec{x} = (\varphi_1(G), \varphi_2(G), \dots, \varphi_n(G))$

2) Dummy player: $V(C \cup \{a_i\}) = V(C) \quad \forall C \subseteq A \setminus \{a_i\}, \varphi_a(G) = 0$

3) Symmetric: $V(C \cup \{a_i\}) = V(C \cup \{a_j\}) \quad \forall C \subseteq A - \{a_i, a_j\}, \varphi_i(G) = \varphi_j(G)$

4) Additive: given $G_1 = (A, v_1)$ and $G_2 = (A, v_2)$, $G = (A, v_1 + v_2)$

$$\varphi_i(G_1) \quad \varphi_i(G_2) \quad \varphi_i(G) = \varphi_i(G_1) + \varphi_i(G_2)$$

$$\vec{x} = (x_1, x_2, \dots, x_n)$$

$$\varphi_1(G) \quad \varphi_2(G)$$

In gran coalition, the only way to satisfying these properties is assign to $x_i, \varphi_i(G)$

In convex games this payoff vector \vec{x} is in the core.

Goodness of CS: $v(CS) = \sum_{c \in CS} v(c)$, $CS \in P^A$ - set of all possible partitions

In super additive games no CS is strictly better than gran coalition, with high # of agents we have a representation problem (2^n subsets, we cannot sum single values).

Dynamic Programming Algorithm: $\varphi(c) = \begin{cases} v(c) & \text{if } |c|=1 \\ \max_{c' \subseteq c} \{v(c'), \max_{c'' \subseteq c \setminus c'} \{\varphi(c') + \varphi(c'')\}\} & \text{otherwise} \end{cases}$

Looking at the best way of splitting a coalition in exactly two sets,

compare the MAX value I obtain with the value of that coalition (remember how I obtained it). Still exponential but ensure to obtain the best solution with reduction of complexity.

Distributed coalition formation - Shohory-Kraus

1. $C_i \leftarrow$ all coalitions including a_i
2. $C_i^* \leftarrow \arg \max_{C \in C_i} v(C)$ same knowledge
for each agent
3. broadcast (a_i, C_i^*) , add received (a_j, C_j^*) to C^*
4. $C_{\max} \leftarrow$ largest C s.t. $a_j \in C$, $(a_j, C) \in C^*$ all the agents in C would like to be part of.
5. if $a_i \in C_{\max}$ then join C_{\max} and return
6. delete from C_i all the coalitions including agents in C_{\max}
7. go to 2.

Not ensured to find the best solution.

MULTIAGENT PLANNING

Calculate a sequence of actions to maximize a global function. In MDP the sequence is implicitly encoded in policy, in STRIPS it is explicitly encoded.

Take in account the long term (a greedy solution could not reach a solution).

Build sequences for each agent, you also have to know what the other agents plan \rightarrow two sources of complexity: length (temporal) and # of agents.

The state is factored: the effect of an action involves a subset of agents \leftarrow algorithms

Ways to intend multiagents:

- centralized algorithm that calculate (multiagent) plan;
- distributed construction of the plan (parallelize it).

Coordination before planning: I decide rules, then every agent builds his plan. The initial coordination ensures that plans are not conflicting (es: driving)

Use Social laws: you assume that every agent knows the rules.

Planning before coordination: two diff plans, then solve conflicts before to start the execution (merge of plans). (es: rearrange a schedule to do the same things of a friend in the meantime).

Decision-theoretic m.a. planning: plan implicitly represented in the policy

Dynamic m.a. planning: conflict ignored before to start the execution, solved at the moment (when discovered)

Assumptions: n agents in a stochastic environment, they are synchronous in performing action until a finite - infinite horizon. Don't know the state of the environment, reward for the system.

DEC-POMDP (Partially observable MDP, decentralized)

$I = \{1, 2, \dots, n\}$ agents

S : states of the world

A_i : action that agent i can perform. $\vec{A} = \bigotimes_{i \in I} A_i$: set of all possible jointed actions (cartesian product)

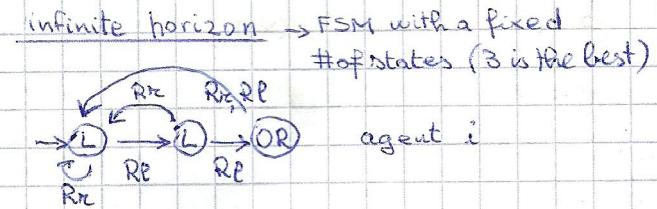
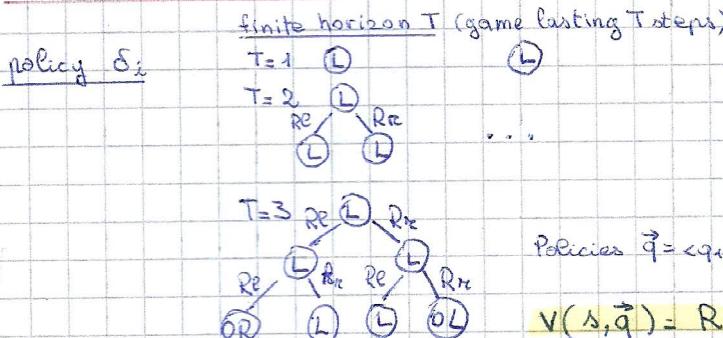
$P: S \times \vec{A} \rightarrow \Delta S$ transition function

Ω_i : observations (set of perceptions of agent i) $\vec{\Omega} = \bigotimes_{i \in I} \Omega_i$

O: $\vec{A} \times S \rightarrow \vec{\Omega}$ observation function: a set of actions brings in a certain state, it gives a probability of making a jointed observation. Independent observations: $O(\vec{a}|s) = O_1(a_1|s) \cdot O_2(a_2|s) \cdot \dots$

R: $S \times \vec{A} \rightarrow \mathbb{R}$ reward function, single function for the system.

Definition of solution is missing (it is in the form of a policy for each agent, we'd like to maximize sth) since the state of the world is not observable, the policy is more complex. Map a sequence of observations to actions. How?



$$\text{Policies } \vec{q} = (q_1, q_2, \dots) \quad V(s_0, \vec{q}) = E \left[\sum_{t=0}^{T-1} R(\vec{a}_t, s_t) \right]$$

$$V(s, \vec{q}) = R(s, \vec{a}) + \sum_{s', \vec{a}'} P(s'|s, \vec{a}) \cdot O(\vec{a}'|s', \vec{a}) \cdot V(s', \vec{q})$$

Find a policy that maximizes $V(s, \vec{q})$ is NEXP-hard. $V(s', \vec{q})$ is used to calculate the result

DCOP (distributed constraint optimization problem)

Every agent decides independently but the goal is to max a global function (exchange info to obtain it)

Constraints network

$X = \{x_1, x_2, \dots\}$ variables

$D = \{D_1, D_2, \dots\}$ domains (finite)

$C = \{C_1, C_2, \dots\}$ constraints: hard $\rightarrow C_i \subseteq D_{x_1} \times D_{x_2}$ (a combination is admissible or not) utility

soft \rightarrow pay a cost if you choose a combination $F_i: D_{x_1} \times D_{x_2} \rightarrow R$

It is possible to write all the constraints as binary constraints, $F_i = -\infty$ simulates an hard constraint

COP: you have a constraint network (represented with constraint graph \rightarrow nodes are the variables, edges are the constraints) and you find an assignment to the variables in the network s.t. it optimizes the sum of soft constraints: $\sum_{F_i \in C} F_i(x_i, x_j)$

DCOP: every agent controls a subset of variables (only one variable). The agent knows the constraints in which the variable is involved and decides the value to assign. NP-hard

Why use DCOP? - scalability: worst-case complexity is the same (in COP big instances are a problem);

- privacy: an agent needs knowledge on domain and constraints about his variable;

- robustness: it works even without an agent

DPOP (dynamic programming optimization problem \rightarrow solving algorithm)

Phase 1 - DFS arrangement (preparation, Depth First Search) \rightarrow build a pseudo-tree starting from a random node and visiting depth first (diff trees for diff choices). I'll add pseudo-arcs to represent all the constraints. Induced width: max # of parents and pseudo-parents a node can have. The smaller it is, the more efficient the algorithm is.

Phase 2 - Util propagation \rightarrow flow of the msgs from the leaf, going up. $U_{i \rightarrow j}(\text{sep}) = \max_{x_i \in X_{k \rightarrow i}} (\sum_{x_k \in C_i} U_{k \rightarrow i} \otimes F_{i,p})$

$U_{i \rightarrow j}$ sent from a child to his real parent.

Sep_i : set of all the agents preceding i in the pseudo-tree who are connected with i or a descendant

C_i : children of ag i P_i, P_{P_i} : parents, pseudo-parents of ag i

Phase 3 - Value propagation \rightarrow flow from the top to the bottom $x_i^* = \arg \max_{x_i \in X_{k \rightarrow i}} (\sum_{k \rightarrow i} U_{k \rightarrow i} + \sum_{j \in P_{P_i}} F_{i,j}(x_i, x_j^*))$

Max size sent is the max size of a separator.

Algorithm optimal but exponential (in the msg size).

$$V_{i \rightarrow j} = \{x_i = x_i^*\} \cup \{x_j = x_j^*\}$$

$\forall s \in Sep_i, s \in Sep_j$

MAX-SUM \rightarrow not guaranteed to provide the optimal solution (could also not converge). Every agent is exchanging iteratively msgs with neighbours. We assume synchronous execution. At every step each agent sends to all the neighbours:

$$m_{i \rightarrow j}^{(t+1)} = \alpha_{i,j} + \max_{x_i} (F_{i,j}(x_i, x_j) + \sum_{k \in N_{i,j}} m_{k \rightarrow i}(x_k))$$

$\downarrow \quad \downarrow \quad \downarrow$
constant for constraint with sum of msgs
receiving agent from neighbor $\neq j$

At the beginning $m_{i \rightarrow j}(x_j) = \emptyset$
If all the agents receive the same msgs two times in a row, the algorithm ends.

$$\sum_{x_j} m_{i \rightarrow j}(x_j) = 0 \text{ for convergence}$$

$$z_i(x_i) = \sum_{k \in N_i} m_{k \rightarrow i}(x_k)$$

$$\hat{x}_i = \arg \max_{x_i} z_i(x_i)$$

If the starting network is acyclic, max-sum converges to the optimal solution. Otherwise, if converges we have a local maximum.

Size of msgs = $|D|$ but you exchange a lot of msgs.

MULTIAGENT LEARNING

Advantages:

- Robustness: if a single agent fails, the others continue to work;
- Efficiency: divide work between agents, the extra cost in coordinate agents is more than compensated;
- Reconfigurable: local changes in the system in dynamic environments;
- Scalable: easy to add new agents;
- Adaptable: you can use it in diff situations → mechanism of learning.

Challenges in learning:

- 1) The space of actions is very large: you have to learn how to connect states in every possible combination of actions.
- 2) Orthogonal problem where the system receives the credits: how to divide credits among the sequence of actions? how to divide among the single agents?

Schema to solve the problems: objective function for the system and local objective functions for the agents.

Factoredness → alignment between local and global objectives (es. if I max locals, I'm max global. Not always);

Learnability → How much my local objective depends on the action of that agent?

Multiagent Reinforcement Learning: max rewards obtained by some actions. Agents don't know the transition function. Learning with trial and error processes: exploration and exploitation.

MARKOV STOCHASTIC GAMES

S: states

$A = A^1 A^2 \dots A^n$ actions # agent

$r = r^1 r^2 \dots r^n$ rewards function $r^i: S \times A^1 \times A^2 \times \dots \times A^n \rightarrow \mathbb{R}$

$p =$ transition function $p: S \times A^1 \times A^2 \times \dots \times A^n \rightarrow \Delta S$

In a synchronous way, the agents perform an action and receive a reward. The world will go to s' .

$$\sum_{s' \in S} p(s'|s, a^1, \dots, a^n) = 1$$

Strategy $\pi^i: S \rightarrow A^i$ Interacting with the world we want to learn sth we don't know (to find π^i).

1. The state of the world encloses the knowledge of all the agents.
2. The policy of the single agent is a path if p is deterministic (in general you have to define an action for every state).

Each agent will max his value function: $v^i(s, \pi^1, \dots, \pi^n) = \sum_{t=0}^{\infty} \beta^t E[\pi_t^i | \pi^1, \pi^2, \dots, \pi^n, s_0=s]$

↓
Problem: find his own strategy to maximize it. Agent-i has no knowledge about the other strategies, it modifies his strategy in order to max v^i , same for the others.

Idea: NASH EQUILIBRIUM $(\pi_1^*, \pi_2^*, \dots, \pi_n^*)$ s.t. $v^i(s, \pi_1^*, \dots, \pi_n^*) \geq v^i(s, \pi_1^*, \dots, \pi_i^{i+1}, \dots, \pi_n^*) \quad \forall s, \forall i, \forall \pi \in \Pi^i$

Nash-Q value: $Q_1^i(s, a^1, a^2, \dots, a^n) = r^i(s, a^1, a^2, \dots, a^n) + \beta \sum_{s' \in S} p(s'|s, a^1, a^2, \dots, a^n) \cdot v^i(s', \pi_1^*, \dots, \pi_n^*)$
The first action is arbitrary, from the next one NE is followed

Stage game: normal-form game played by two agents.

Once you have it, you calculate the NE of the game.

What if agents don't know r, p ? They try, knowledge from what happens. Agent perceives the state of the environment, take an action and see:

- 1) reward I get
- 2) reward others get
- 3) actions that other agents took
- 4) new state

$$Q_{t+1}^i(s, a^1, \dots, a^n) = (1-\alpha_t) Q_t^i(s, a^1, \dots, a^n) + \alpha_t [\pi_t^i + \beta \text{Nash } Q_t^i(s')]$$

$\alpha_t \rightarrow$ new value = old value $\alpha_t = 1 \rightarrow$ completely forget the old value. Every agent keeps updated all Q^i .

Nash $Q_t^i(s')$: payoff that ag*i* gets in the stage game of s' playing a NE. All the agents choose the same NE between multiples.

↳ have to build the stage game. $Q_{t=0}^i(\dots) = \emptyset$

$\frac{1}{n} \sum_{t=1}^n Q_t^i(s, a^1, \dots, a^n) \rightarrow$ # of times I saw this combination: exploration at the beginning, exploitation later.
Idea: continuing to update, $Q_t^i \rightarrow Q_*^i$ (build what you don't know: r, p).
↓ called

Nash-Q learning: it brings me a NE exploring a huge space.

a_1	a_2	a_3
a_1	$Q_1^*(s, a_1, a_2)$	$Q_1^*(s, a_1, a_3)$
a_2	$Q_2^*(s, a_1, a_2)$	$Q_2^*(s, a_1, a_3)$
a_3	$Q_3^*(s, a_2, a_1)$	$Q_3^*(s, a_2, a_3)$
a_1	$Q_1^*(s, a_3, a_2)$	$Q_1^*(s, a_3, a_1)$