Lifting in Multi-agent Systems under Uncertainty

Tanya Braun¹, Marcel Gehrke², Florian Lau², Ralf Möller²

 $^1 \rm Computer$ Science Department, University of Münster, $^2 \rm Institute$ of Information Systems, University of Lübeck, $^3 \rm Institute$ of Telematics, University of Lübeck

Multi-agent Systems

- Joint state and reward for transition, sensor, reward functions
- Individual actions and observations
 Complexity exponential in number of agents n
 ↓
 Problems with large n not computable, e.g., nano-scale systems

Under Symmetries

- Agent types: same actions and observations available

Lifting in Multi-agent Systems under Uncertainty

Tanya Braun¹, Marcel Gehrke², Florian Lau², Ralf Möller²

¹University of Münster, Münster, Germany ²University of Lübeck, Lübeck, Germany

DecPOMDP

- Decentralised Partially Observable Markov Decision Problem
- Set of agents working towards a joint reward
- Environment a *Markov-1 stochastic process*
- DecPOMDP model M

 $(I, S, \{A_i\}_{i=1}^N, T, R, \{O_i\}_{i=1}^N, \Omega)$

- I a set of agents, |I| = N,
- S a random variable for the state space,
- A_i a decision variable; $\mathbf{A} = \times_{i=1}^N ran(A_i)$ set of joint actions
- T(S', S, A) = P(S' | S, A) a transition model,
- R(S) a reward function,
- O_i a random variable; $O = \times_{i=1}^N ran(O_i)$ set of joint observations,
- $\Omega(\mathbf{0}, S) = P(\mathbf{0} \mid S)$ a sensor model.
- Given horizon τ , each agent *i* has a local policy π_i : $ran(O_{i,(0:t)}) \mapsto ran(A_i); t < \tau; \boldsymbol{\pi} = (\pi_i)_{i=1}^N$ a joint policy
- DecPOMDP Semantics: all possible joint policies Π_M
- *DecPOMDP Problem*: find joint policy π^* that maximises the expected utility with discount factor $\gamma \in [0,1]$

$$U_{M}^{\pi}(s_{t}, \boldsymbol{o}_{0:t}) = R(s_{t}, \frac{\pi(\boldsymbol{o}_{0:t})}{a_{0:t}}) + \gamma^{t} \sum_{s_{t+1} \in ran(S)} T(s_{t+1}, s_{t}, \pi(\boldsymbol{o}_{0:t}))$$
$$\cdot \sum_{\boldsymbol{o}_{t+1} \in ran(\boldsymbol{o})} \Omega(\boldsymbol{o}_{t+1}, s_{t+1}) U_{M}^{\pi}(s_{t+1}, \boldsymbol{o}_{0:t+1})$$

Problem

• **Complexities** \mathbb{T} , \mathbb{R} , \mathbb{O} of T, R, Ω , evaluation cost \mathbb{C} of a joint policy and size of policy space \mathbb{P} *exponential* in N

 $\mathbb{T} \in O(s^2 a^N) \qquad \mathbb{R} \in O(sa^N) \qquad \mathbb{O} \in O(so^N)$ $\mathbb{C} \in O(so^{N\tau}) \qquad \mathbb{P} \in O(a^{\frac{N(o^{\tau}-1)}{o-1}})$

• $s = |ran(S)|, a = \max_{i} |ran(A_i)|, o = \max_{i} |ran(O_i)|$

 \rightarrow Especially a problem with $N \gg 10,000$

Symmetric & Partitioned DecPOMDP

- Symmetric DecPOMDP: *K* << *N* partitions in agent set
- In each partition: the same set of actions and observations
- Formally, $I = \bigcup_{k=1}^{K} \mathfrak{I}_k$, $\mathfrak{I}_k \neq \emptyset$, $\mathfrak{I}_k \cap \mathfrak{I}_{k'} = \emptyset$; for each \mathfrak{I}_k : $\forall i, j \in \mathfrak{I}_k : ran(A_i) = ran(A_i) \land ran(O_i) = ran(O_i)$
- *Partitioned* model
 - $(\overline{I}, S, \{\overline{A}_k\}_{k=1}^K, \overline{T}, \overline{R}, \{\overline{O}_k\}_{k=1}^K, \overline{\Omega})$
- \overline{I} a partitioning $\{\mathfrak{I}_k\}_{k=1}^K$ of an agent set, $n_k = |\mathfrak{I}_k|, |\overline{I}| = N$,
- \overline{A}_k , \overline{O}_k variables per partition, with joint actions and observations defined over the *K* partitions
- $\overline{T}, \overline{R}, \overline{\Omega}$ like T, R, Ω but defined over partitioned inputs

Counting DecPOMDP

$$A_k = \overline{O}_k = \overline{O}_k = \overline{T}, \overline{R}$$

- Theorem:
- Theorem:

Optimal policy in M_c also optimal in $gr(M_c)$. • Complexities \mathbb{T}_c , \mathbb{R}_c , \mathbb{O}_c , \mathbb{C}_c *polynomial*, \mathbb{P}_c *exponential* in n < N $\mathbb{T}_c \in O(s^2 n^{Ka}) \qquad \mathbb{R}_c \in O(s n^{Ka}) \qquad \mathbb{O}_c \in O(s n^{Ko})$ $\mathbb{C}_{c} \in O(sn^{K\tau o}) \qquad \mathbb{P}_{c} \in O(n^{a^{\frac{K(n^{\tau o}-1)}{n^{o}-1}}})$ • $n = \max_{k} n_{k}, |ran(\bar{A}_{k})| \le n^{a}, |ran(\bar{O}_{k})| \le n^{o}$

Isomorphic DecPOMDP

- Lemma:
- Lemma:
- Corollary:

Indistinguishability of agents in a partition yields invariance towards which particular partition agents perform an action or observe an event: it only matters how many agents do or observe something • E.g., 2 actions performed by 5 agents each: 10 over 5 and 5 = 255 different permutations of 10 agents to do the actions with the same outcome

 \rightarrow Count occurrences $[n_1, \dots, n_l], l = |ran(A_k)|$

• *Counting DecPOMDP*: a partitioned DecPOMDP with

 $= #_{X_k}[A_k(X_k)]$ a counting random variable

 $= #_{X_k}[O_k(X_k)]$ a counting random variable

 $\overline{\Omega}$, $\overline{\Omega}$ defined over the counting random variables

 $\left(\overline{I}, S, \left\{\#_{X_k}[A_k(X_k)]\right\}_{k=1}^K, \overline{T}, \overline{R}, \left\{\#_{X_k}[O_k(X_k)]\right\}_{k=1}^K, \overline{\Omega}\right)$

Counting model M_c has an equivalent ground model $gr(M_c)$.

• Independence assumption between agents of a partition $\rightarrow \overline{T}, \overline{R}, \overline{\Omega}$ factorise identically for each agent within each partition \rightarrow Higher efficiency at the expense of lower expressiveness • *Isomorphic DecPOMDP*: a partitioned DecPOMDP with $\bar{A}_k = A_k(X_k)$ a parameterised random variable • $\overline{O}_k = O_k(X_k)$ a parameterised random variable • \overline{T} , \overline{R} , $\overline{\Omega}$ defined over the parameterised random variables $(\overline{I}, S, \{A_k(X_k)\}_{k=1}^K, \overline{T}, \overline{R}, \{O_k(X_k)\}_{k=1}^K, \overline{\Omega})$

Isomorphic model M_i has an equivalent counting model M_c .

Enough to define policies in M_i over $ran(A_k)$ and $ran(O_k)$.

Partition sizes only influence U_M^{π} , not π^* itself. • Complexities \mathbb{T}_i , \mathbb{R}_i , \mathbb{O}_i , \mathbb{C}_i , \mathbb{P}_i *logarithmic* in n < N $\mathbb{T}_i \in O(s^2 a^K) \qquad \mathbb{R}_i \in O(s a^K) \qquad \mathbb{O}_i \in O(s o^K)$ $\mathbb{C}_i \in O(\log_2(n) s o^{K\tau}) \qquad \mathbb{P}_i \in O(a^{\frac{K(o^{\tau}-1)}{o-1}})$

Possible to reuse existing solution approaches to solve DecPOMDP for K partitions; adapt result w.r.t. sizes n_k to reach a given U

Example: DecTiger

- $I = \{agent_1, agent_2\}, \Im_{K=1} = I$
- $ran(S) = \{tigerleft, tigerright\}$
- $A_K(X_K) = \{ listen, openleft, openright \}$
- $O_K(X_K) = \{hearleft, hearright\}$
- As a counting model: T, R, Ω equivalently encoded in $\overline{T}, \overline{R}, \overline{\Omega}$
- As an isomorphic model:
- Reset in T not possible to encode in \overline{T}
- Agreeing on an action yielding a lower punishment not possible to encode in \overline{R}
- Ω equivalently encoded in $\overline{\Omega}$

Application: Nanoscale Medical Systems







ink to preliminary paper (QR code): https://www.ifis.uni-luebeck.de/upl oads/tx wapublications/braun 143 public.pdf Please cite full version when published by PMLR.



Agent set: $\kappa = 4$ different types of markers/sensors, $\iota = 1$ different types of messages/nanobots; $\kappa + \iota = K$, $n_k \sim 64,000$ (preliminary experiments) Each \mathfrak{I}_k : 2 actions (output, no act.), 2 observations (sense/receive, no obs.)

