

Dynamic programming in infinite horizon dynamic teams with non-classical information structure

Aditya Mahajan

McGill University

Fifth Workshop on Dynamic Games in Management Science, GERAD, Montreal
28 Nov, 2014

Joint work with: Ashutosh Nayyar and Demosthenis Teneketzis

What are dynamic teams

Dynamic games where all players have identical payoff function

What are dynamic teams

Dynamic games where all players have identical payoff function

Literature ▶ Economics Literature

overview

- ▶ Radner, "Team decision problems," Ann Math Stat, 1962.
- ▶ Marschak and Radner, "Economics Theory of Teams," 1972.
- ▶ ...

▶ Systems & Control Literature

- ▶ Witsenhausen, "Separation of estimation and control," Proc IEEE, 1971.
- ▶ Witsenhausen, "On information structures, feedback and causality," SICON 1971.
- ▶ Ho and Chu, "Team decision theory and information structures," IEEE TAC 1972.
- ▶ ...

What are dynamic teams

Dynamic games where all players have identical payoff function

Literature ▶ Economics Literature

overview

- ▶ Radner, "Team decision problems," Ann Math Stat, 1962.
- ▶ Marschak and Radner, "Economics Theory of Teams," 1972.
- ▶ ...

▶ Systems & Control Literature

- ▶ Witsenhausen, "Separation of estimation and control," Proc IEEE, 1971.
- ▶ Witsenhausen, "On information structures, feedback and causality," SICON 1971.
- ▶ Ho and Chu, "Team decision theory and information structures," IEEE TAC 1972.
- ▶ ...

Simpler than non-cooperative game theory.

All "pre-game" agreements are enforceable.

Simpler than cooperative game theory.

The value of the game does not need to be split between the players.

What are dynamic teams

Dynamic games where all players have identical payoff function

Main difficulties:

- ▶ Information decentralization
Non-classical information structure
- ▶ Seeking global optimality

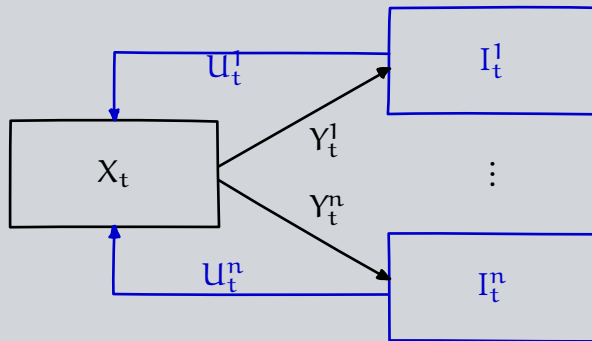
Simpl

All “pre-game” agreements are enforceable.

Simpler than cooperative game theory.

The value of the game does not need to be split between the players.

Simplest general model of a dynamic team



Dynamics $X_{t+1} = f_t(X_t, \mathbf{u}_t, W_t^0)$, where $\mathbf{u}_t = (u_t^1, \dots, u_t^n)$.

Observation $Y_t^i = h_t^i(X_t, W_t^i)$.

Information structure

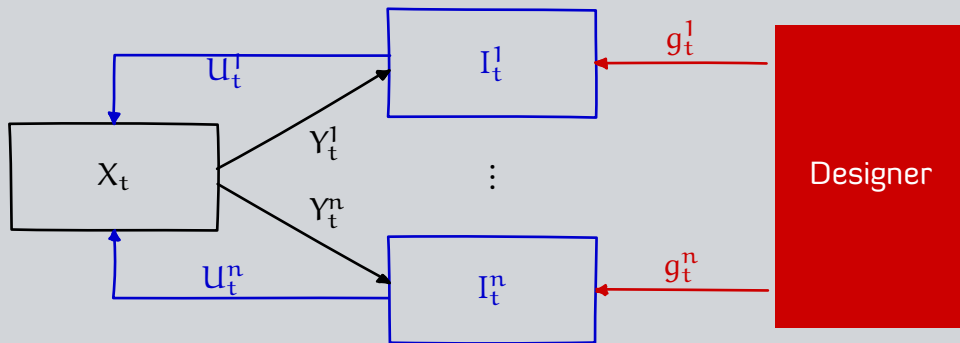
$$\{Y_{1:t}^i, u_{1:t-1}^i\} \subseteq \mathbf{I}_t^i \subseteq \{\mathbf{Y}_{1:t}, \mathbf{u}_{1:t-1}\}, \quad u_t^i = g_t^i(I_t^i).$$

Control Strategy $\mathbf{g} = (g^1, \dots, g^n)$, where $g^i = (g_1^i, g_2^i, \dots)$.

Performance ▶ Per-step reward $R_t = \rho(X_t, \mathbf{u}_t)$. ▶

$$J(\mathbf{g}) = \mathbb{E}^g \left[\sum_{t=0}^{\infty} \beta^t R_t \right]$$

Simplest general model of a dynamic team



Dynamics $X_{t+1} = f_t(X_t, \mathbf{U}_t, W_t^0)$, where $\mathbf{U}_t = (U_t^1, \dots, U_t^n)$.

Observation $Y_t^i = h_t^i(X_t, W_t^i)$.

Information structure

$$\{Y_{1:t}^i, U_{1:t-1}^i\} \subseteq \mathbf{I}_t^i \subseteq \{Y_{1:t}, \mathbf{U}_{1:t-1}\}, \quad U_t^i = g_t^i(I_t^i).$$

Control Strategy $\mathbf{g} = (g^1, \dots, g^n)$, where $g^i = (g_1^i, g_2^i, \dots)$.

Performance ▶ Per-step reward $R_t = \rho(X_t, \mathbf{U}_t)$. ▶

$$J(\mathbf{g}) = \mathbb{E}^{\mathbf{g}} \left[\sum_{t=0}^{\infty} \beta^t R_t \right]$$

Conceptual difficulties

The optimal control problem is a **functional optimization** problem where we have to choose an **infinite sequence of control laws g** to maximize the expected total reward.

The domain I_t^i of control law g_t^i increases with time.

- ▶ Can the optimization problem be solved?
- ▶ Can we implement the optimal solution?

Agent based methods lead to infinite regress.

Signaling (or the communication aspect of control)

Centralized stochastic control: Information state

$$I_t \subseteq I_{t+1}$$

Centralized stochastic control: Information state

$$I_t \subseteq I_{t+1}$$

A process $\{Z_t\}_{t=0}^{\infty}$ is called an **information state** if

► **Function of available information**

There exists a series of functions $\{F_t\}_{t=0}^{\infty}$ such that $Z_t = f_t(I_t)$.

► **Absorbs the effect of available information on current rewards**

$$\mathbb{P}(R_t \in \mathcal{B} \mid I_t = i_t, U_t = u_t) = \mathbb{P}(R_t \in \mathcal{B} \mid Z_t = F_t(i_t), U_t = u_t).$$

► **Controlled Markov property**

$$\mathbb{P}(Z_{t+1} \in \mathcal{A} \mid I_t = i_t, U_t = u_t) = \mathbb{P}(Z_{t+1} \in \mathcal{A} \mid Z_t = F_t(i_t), U_t = u_t).$$

Examples: ► System state in MDPs ► Belief state in POMDPs

Centralized control: Structure of optimal strategies

The information state absorbs the effect of available information on **expected future cost**, i.e., for any choice of **future strategy** $\mathbf{g}_{(t)} = (g_{t+1}, g_{t+2}, \dots)$

$$\mathbb{E}^{\mathbf{g}_{(t)}} \left[\sum_{\tau=t}^{\infty} \beta^{\tau} R_{\tau} \mid I_t = i_t, U_t = u_t \right] = \mathbb{E}^{\mathbf{g}_{(t)}} \left[\sum_{\tau=t}^{\infty} \beta^{\tau} R_{\tau} \mid Z_t = F_t(i_t), U_t = u_t \right].$$

Centralized control: Structure of optimal strategies

The information state absorbs the effect of available information on **expected future cost**, i.e., for any choice of **future strategy** $\mathbf{g}_{(t)} = (g_{t+1}, g_{t+2}, \dots)$

$$\mathbb{E}^{\mathbf{g}_{(t)}} \left[\sum_{\tau=t}^{\infty} \beta^{\tau} R_{\tau} \mid I_t = i_t, U_t = u_t \right] = \mathbb{E}^{\mathbf{g}_{(t)}} \left[\sum_{\tau=t}^{\infty} \beta^{\tau} R_{\tau} \mid Z_t = F_t(i_t), U_t = u_t \right].$$

Therefore,

- ▶ Z_t is a sufficient statistic for performance evaluation,
- ▶ there is **no loss of optimality** in using control laws of the form $g_t: Z_t \mapsto U_t$

Centralized control: Structure of optimal strategies

The information state absorbs the effect of available information on **expected future cost**, i.e., for any choice of **future strategy** $\mathbf{g}_{(t)} = (g_{t+1}, g_{t+2}, \dots)$

$$\mathbb{E}^{\mathbf{g}_{(t)}} \left[\sum_{\tau=t}^{\infty} \beta^{\tau} R_{\tau} \mid I_t = i_t, U_t = u_t \right] = \mathbb{E}^{\mathbf{g}_{(t)}} \left[\sum_{\tau=t}^{\infty} \beta^{\tau} R_{\tau} \mid Z_t = F_t(i_t), U_t = u_t \right].$$

Therefore,

- ▶ Z_t is a sufficient statistic for performance evaluation,
- ▶ there is **no loss of optimality** in using control laws of the form $g_t: Z_t \mapsto U_t$

- Examples**
- ▶ In MDPs, $g_t: X_t \mapsto U_t$.
 - ▶ In POMDPs, $g_t: B_t \mapsto U_t$, where B_t is the belief state.

Centralized control: Dynamic programming

For any strategy \mathbf{g} of the form $g_t: Z_t \mapsto U_t$,

$$\begin{aligned} \mathbb{E}^{\mathbf{g}^{(t)}} \left[\mathbb{E}^{\mathbf{g}^{(t+1)}} \left[\sum_{\tau=t+1}^{\infty} \beta^{\tau} R_{\tau} \mid Z_{t+1}, U_{t+1} = g_{t+1}(Z_{t+1}) \right] \mid Z_t = z_t, U_t = u_t \right] \\ = \mathbb{E}^{\mathbf{g}^{(t)}} \left[\sum_{\tau=t+1}^{\infty} \beta^{\tau} R_{\tau} \mid Z_t = z_t, U_t = u_t \right] \end{aligned}$$

Relies on $I_t \subseteq I_{t+1}$

Centralized control: Dynamic programming

For any strategy g of the form $g_t: Z_t \mapsto U_t$,

$$\begin{aligned} \mathbb{E}^{g^{(t)}} \left[\mathbb{E}^{g^{(t+1)}} \left[\sum_{\tau=t+1}^{\infty} \beta^{\tau} R_{\tau} \mid Z_{t+1}, U_{t+1} = g_{t+1}(Z_{t+1}) \right] \mid Z_t = z_t, U_t = u_t \right] \\ = \mathbb{E}^{g^{(t)}} \left[\sum_{\tau=t+1}^{\infty} \beta^{\tau} R_{\tau} \mid Z_t = z_t, U_t = u_t \right] \quad \text{Relies on } I_t \subseteq I_{t+1} \end{aligned}$$

There exists a **time-homogeneous** optimal strategy $g^* = (g^*, g^*, \dots)$ that is given by the fixed point of the following dynamic program

$$V(z) = \min_{u \in U} \mathbb{E}[R_t + \beta V(Z_{t+1}) \mid Z_t = z, U_t = u]$$

Centralized control: Dynamic programming

For any strategy g of the form $g_t: Z_t \mapsto U_t$,

$$\mathbb{E}^{g_0} \left[\sum_{t=0}^{\infty} \gamma^t U_t \mid Z_0 = z_0 \right]$$

Both these results rely on an appropriate choice of
information state.

Note that information state for DP
is also a sufficient statistic for control.

Centralized control: Dynamic programming

For any strategy g of the form $g_t: Z_t \mapsto U_t$,

$$\mathbb{E}^g \left[\sum_{t=0}^{\infty} \gamma^t U_t \mid Z_0 = z_0 \right]$$

- ▶ Can we identify a **sufficient statistic** Z_t^i and restrict attention to $g_t^i: Z_t^i \mapsto U_t^i$?
- ▶ Can we show that there exist **time-homogeneous** optimal control strategies?
- ▶ Can we identify appropriate **information states** to determine a **dynamic program** that computes such optimal strategies?

Two approaches to dynamic programming:
The person-by-person approach

The person-by-person approach

Pick an agent, say i .

Arbitrarily fix the strategies g^{-i} of all other agents.

Identify an information-state process $\{Z_t^i\}_{t=0}^\infty$ for agent i .

Structure of optimal strategies If \mathcal{Z}_t^i , the space of realization of Z_t^i , does not depend on g^{-i} , then there is no loss of optimality in using $g_t^i: \mathcal{Z}_t^i \mapsto \mathcal{U}_t^i$.

-
- ▶ Radner, "Team decision problems," Ann Math Stat, 1962.
 - ▶ Marschak and Radner, "Economics Theory of Teams," 1972.

The person-by-person approach

Pick an agent, say i .

Arbitrarily fix the strategies g^{-i} of all other agents.

Identify an information-state process $\{Z_t^i\}_{t=0}^\infty$ for agent i .

Structure of optimal strategies If \mathcal{Z}_t^i , the space of realization of Z_t^i , does not depend on g^{-i} , then there is no loss of optimality in using $g_t^i: \mathcal{Z}_t^i \mapsto \mathcal{U}_t^i$.

Write coupled dynamic programs to identify the best response strategy

$$g^i = \mathcal{D}^i(g^{-i})$$

- Remarks**
- ▶ Is the best-response strategy time-homogeneous?
 - ▶ The coupled dynamic program always has a fixed-point.
 - ▶ Is the fixed point unique?

-
- ▶ Radner, "Team decision problems," Ann Math Stat, 1962.
 - ▶ Marschak and Radner, "Economics Theory of Teams," 1972.

The person-by-person approach

Pick an agent, say i .

The person-by-person approach:

- ▶ May identify the **structure** of globally optimal control strategies.
- ▶ Provides coupled dynamic programs, which, at best, may determine **person-by-person** optimal control strategies. Such strategies can be **arbitrarily bad** compared to globally optimal strategies.

Remarks

- ▶ Is the best-response strategy **time-homogeneous**?
- ▶ The coupled dynamic program always has a fixed-point.
- ▶ Is the fixed point unique?

- ▶ Radner, "Team decision problems," Ann Math Stat, 1962.
- ▶ Marschak and Radner, "Economics Theory of Teams," 1972.

An example: coupled subsystems with control sharing

Global state $\mathbf{X}_t = (X_t^1, \dots, X_t^n)$

Dynamics $X_{t+1}^i = f^i(X_t^i, \mathbf{u}_t, W_t^i)$, where $\mathbf{u}_t = (u_t^1, \dots, u_t^n)$.

Information
structure

$$I_t^i = \{X_{1:t}^i, \mathbf{u}_{1:t-1}\}$$

An example: coupled subsystems with control sharing

Global state $\mathbf{X}_t = (X_t^1, \dots, X_t^n)$

Dynamics $X_{t+1}^i = f^i(X_t^i, \mathbf{u}_t, W_t^i), \quad \text{where } \mathbf{u}_t = (u_t^1, \dots, u_t^n).$

Information
structure

$$I_t^i = \{X_{1:t}^i, \mathbf{u}_{1:t-1}\}$$

Conditional
independence

For any arbitrary choice of control strategies \mathbf{g} :

$$\mathbb{P}(\mathbf{X}_{1:t} \mid \mathbf{u}_{1:t-1} = \mathbf{u}_{1:t-1}) = \prod_{i=1}^n \mathbb{P}(X_{1:t}^i \mid \mathbf{u}_{1:t-1} = \mathbf{u}_{1:t-1})$$

An example: coupled subsystems with control sharing

Global state $\mathbf{X}_t = (X_t^1, \dots, X_t^n)$

Dynamics $X_{t+1}^i = f^i(X_t^i, \mathbf{u}_t, W_t^i)$, where $\mathbf{u}_t = (u_t^1, \dots, u_t^n)$.

Information
structure

$$I_t^i = \{X_{1:t}^i, \mathbf{u}_{1:t-1}\}$$

Conditional
independence

For any arbitrary choice of control strategies \mathbf{g} :

$$\mathbb{P}(\mathbf{X}_{1:t} \mid \mathbf{u}_{1:t-1} = \mathbf{u}_{1:t-1}) = \prod_{i=1}^n \mathbb{P}(X_{1:t}^i \mid \mathbf{u}_{1:t-1} = \mathbf{u}_{1:t-1})$$

Structure
of optimal
strategies

- ▶ Arbitrarily fix strategies \mathbf{g}^{-i} , and consider the “best-response” strategy at agent i .
- ▶ $\{X_t^i, \mathbf{u}_{1:t-1}\}$ is an information-state at agent i .

Two approaches to dynamic programming:
The common-information approach

Illustrative example: Two agent static Bayesian team

Observations $I^1 = (C, Y^1)$ and $I^2 = (C, Y^2)$. Joint probability on (X, Y^1, Y^2, C) .

Actions $U^1 = g^1(C, Y^1)$ and $U^2 = g^2(C, Y^2)$.

Utility Maximize $J(g^1, g^2) = \mathbb{E}^{(g^1, g^2)}[\rho(X, U^1, U^2)]$.

Illustrative example: Two agent static Bayesian team

Observations $I^1 = (C, Y^1)$ and $I^2 = (C, Y^2)$. Joint probability on (X, Y^1, Y^2, C) .

Actions $U^1 = g^1(C, Y^1)$ and $U^2 = g^2(C, Y^2)$.

Utility Maximize $J(g^1, g^2) = \mathbb{E}^{(g^1, g^2)}[\rho(X, U^1, U^2)]$.

Complexity Possibilities for $g^i = |\mathcal{U}^i|^{|\mathcal{C}| \cdot |Y^i|}$

Illustrative example: Two agent static Bayesian team

Observations $I^1 = (C, Y^1)$ and $I^2 = (C, Y^2)$. Joint probability on (X, Y^1, Y^2, C) .

Actions $U^1 = g^1(C, Y^1)$ and $U^2 = g^2(C, Y^2)$.

Utility Maximize $J(g^1, g^2) = \mathbb{E}^{(g^1, g^2)}[\rho(X, U^1, U^2)]$.

Complexity Possibilities for $g^i = |\mathcal{U}^i|^{|\mathcal{C}| \cdot |\mathcal{Y}^i|}$

Partial strategies

- ▶ Define a functional-valued function ψ^i such that $g^i(C, Y^i) = \psi^i(C)(Y^i)$.
- ▶ We call $\gamma_C^i = \psi^i(C)$ a **partial strategy** (or prescriptions)

Illustrative example: Two agent static Bayesian team

Observations $I^1 = (C, Y^1)$ and $I^2 = (C, Y^2)$. Joint probability on (X, Y^1, Y^2, C) .

Actions $U^1 = g^1(C, Y^1)$ and $U^2 = g^2(C, Y^2)$.

Utility Maximize $J(g^1, g^2) = \mathbb{E}^{(g^1, g^2)}[\rho(X, U^1, U^2)]$.

Complexity Possibilities for $g^i = |\mathcal{U}^i|^{|\mathcal{C}| \cdot |\mathcal{Y}^i|}$

Partial strategies

- ▶ Define a functional-valued function ψ^i such that $g^i(C, Y^i) = \psi^i(C)(Y^i)$.
- ▶ We call $\gamma_C^i = \psi^i(C)$ a **partial strategy** (or prescriptions)

The idea of a coordinator

Consider a **virtual** coordinator that:

- ▶ Observes the **common information** C
- ▶ And prescribes the partial strategy (γ_C^1, γ_C^2)
- ▶ The decision rule at the coordinator is $\psi: C \mapsto (\gamma_C^1, \gamma_C^2)$.

$$\tilde{J}(\psi) = \mathbb{E}^\psi[\rho(X, \gamma_C^1(Y^1), \gamma_C^2(Y^2))]$$

Illustrative example: Two agent static Bayesian team

- ▶ Centralized optimization problem

Can be solved by solving looking at the conditional reward

$$\mathbb{E}[\rho(X, \gamma_C^1(Y^1), \gamma_C^2(Y^2) \mid C = c)].$$

- ▶ Complexity: $|C| \cdot |\mathcal{U}^1|^{|\mathcal{Y}^1|} \cdot |\mathcal{U}^2|^{|\mathcal{Y}^2|}$

$(Y^i).$

- ▶ Contrast from: $|\mathcal{U}^1|^{|\mathcal{Y}^1|} \cdot |C| \cdot |\mathcal{U}^2|^{|\mathcal{Y}^2|} \cdot |C|$

▶ The decision rule at the coordinator is $\psi: C \mapsto \gamma_C^1(Y_C^1, Y_C^2)$.

$$\tilde{J}(\psi) = \mathbb{E}^\psi[\rho(X, \gamma_C^1(Y^1), \gamma_C^2(Y^2))]$$

One dynamic program to rule them all

$$V(\blacksquare) = \min_{\blacksquare} \mathbb{E}[R_t + \beta V(\blacksquare_{t+1}) \mid \blacksquare_t = \blacksquare, \blacksquare_t = \blacksquare]$$

One dynamic program to rule them all

$$V(\blacksquare) = \min_{\blacksquare} \mathbb{E}[R_t + \beta V(\blacksquare_{t+1}) \mid \blacksquare_t = \blacksquare, \blacksquare_t = \blacksquare]$$

- ▶ The **information state** must be a function of the information available to **every** player.

One dynamic program to rule them all

$$V(\blacksquare) = \min_{\blacksquare} \mathbb{E}[R_t + \beta V(\blacksquare_{t+1}) \mid \blacksquare_t = \blacksquare, \blacksquare_t = \blacksquare]$$

- ▶ The **information state** must be a function of the information available to **every** player.

$$\text{Common information: } C_t = \bigcap_{\tau \geq t} \bigcap_{i=1}^n I_{\tau}^i, \quad \text{Local information: } L_t^i = I_t^i \setminus C_t$$

One dynamic program to rule them all

$$V(z) = \min_{\blacksquare} \mathbb{E}[R_t + \beta V(Z_{t+1}) \mid Z_t = z, \blacksquare_t = \blacksquare]$$

- ▶ The **information state** must be a function of the information available to **every** player.

$$\text{Common information: } C_t = \bigcap_{\tau \geq t} \bigcap_{i=1}^n I_{\tau}^i, \quad \text{Local information: } L_t^i = I_t^i \setminus C_t$$

One dynamic program to rule them all

$$V(z) = \min_{\blacksquare} \mathbb{E}[R_t + \beta V(Z_{t+1}) \mid Z_t = z, \blacksquare_t = \blacksquare]$$

- ▶ The **information state** must be a function of the information available to **every** player.

$$\text{Common information: } C_t = \bigcap_{\tau \geq t} \bigcap_{i=1}^n I_{\tau}^i, \quad \text{Local information: } L_t^i = I_t^i \setminus C_t$$

- ▶ Each step of the dynamic programming must determine a mapping from $(C_t, L_t^i) \mapsto U_t^i$.
 - ▶ The information state Z_t only depends on C_t
 - ▶ Thus, the “action” at each step must be a mapping $L_t^i \mapsto U_t^i$. Call it **prescription** and denote it by γ_t^i .

One dynamic program to rule them all

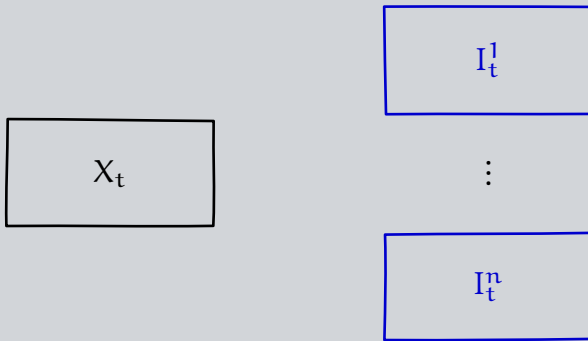
$$V(z) = \min_{\gamma} \mathbb{E}[R_t + \beta V(Z_{t+1}) \mid Z_t = z, \Gamma_t = \gamma]$$

- ▶ The **information state** must be a function of the information available to **every** player.

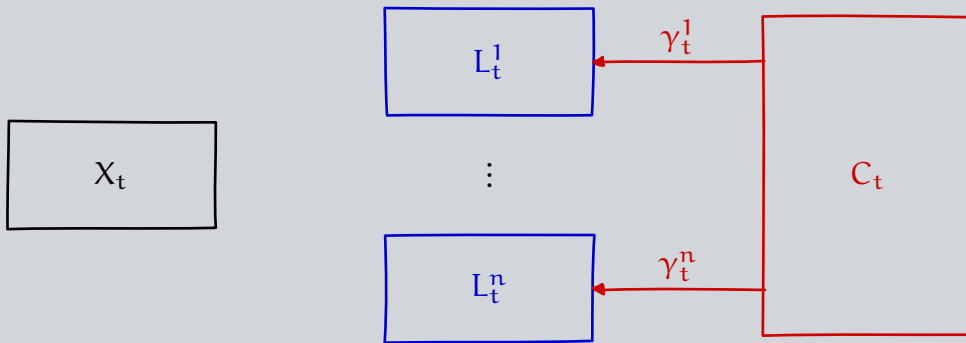
$$\text{Common information: } C_t = \bigcap_{\tau \geq t} \bigcap_{i=1}^n I_{\tau}^i, \quad \text{Local information: } L_t^i = I_t^i \setminus C_t$$

- ▶ Each step of the dynamic programming must determine a mapping from $(C_t, L_t^i) \mapsto U_t^i$.
 - ▶ The information state Z_t only depends on C_t
 - ▶ Thus, the “action” at each step must be a mapping $L_t^i \mapsto U_t^i$. Call it **prescription** and denote it by γ_t^i .

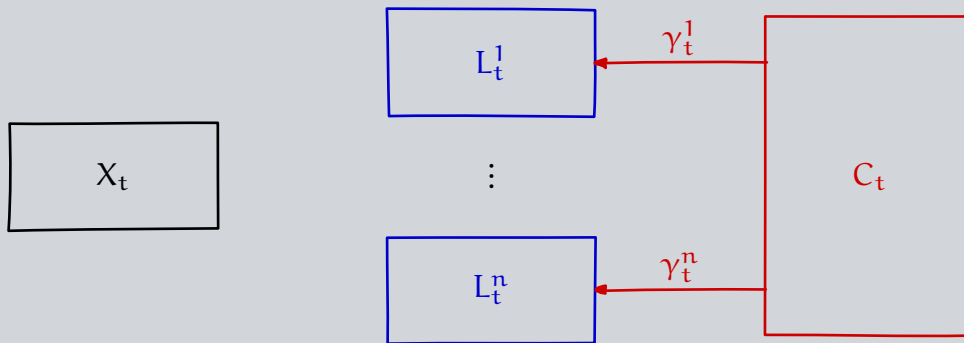
A virtual coordinator



A virtual coordinator



A virtual coordinator



When does this work: Partial history sharing

- ▶ $|\mathcal{L}_t^i|$ is **uniformly bounded (over i and t)** and

$$\mathbb{P}(L_{t+1}^i \in \mathcal{A} \mid \mathbf{C}_t, L_t^i, U_t^i, Y_{t+1}^i) = \mathbb{P}(L_{t+1}^i \in \mathcal{A} \mid L_t^i, U_t^i, Y_{t+1}^i)$$

Centralized POMDP

- ▶ Information state: $\mathbb{P}(X_t, \mathbf{L}_t \mid \mathbf{C}_t = \mathbf{c})$ (or something else)
- ▶ “Standard” POMDP results apply, value function is PWLC.
- ▶ Subsumes many previous results on DP for decentralized stochastic control.

Example 1: Delayed sharing information structure

Dynamics $X_{t+1} = f_t(X_t, \mathbf{U}_t, W_t^0)$, where $\mathbf{U}_t = (U_t^1, \dots, U_t^n)$.

Observations $Y_t^i = h_t^i(X_t, W_t^i)$.

Information structure $I_t^i = \{Y_{1:t}^i, U_{1:t-1}^i, \mathbf{Y}_{1:t-k}, \mathbf{U}_{1:t-k}\}$. k is the sharing delay.

-
- ▶ Witsenhausen, "Separation of estimation and control," Proc IEEE, 1971.
 - ▶ Nayyar, Mahajan and Teneketzis, "Optimal control strategies in delayed sharing information structures," IEEE TAC 2011.

Example 1: Delayed sharing information structure

Dynamics $X_{t+1} = f_t(X_t, \mathbf{U}_t, W_t^0)$, where $\mathbf{U}_t = (U_t^1, \dots, U_t^n)$.

Observations $Y_t^i = h_t^i(X_t, W_t^i)$.

Information structure $I_t^i = \{Y_{1:t}^i, U_{1:t-1}^i, \mathbf{Y}_{1:t-k}, \mathbf{U}_{1:t-k}\}$. k is the sharing delay.

Common info.: $C_t = \{\mathbf{Y}_{1:t-k}, \mathbf{U}_{1:t-k}\}$, Local Info.: $L_t^i = I_t^i \setminus C_t$, Pres.: $\Gamma_t^i: L_t^i \mapsto U_t^i$

Information State $\Pi_t = \mathbb{P}(X_t, \mathbf{L}_t \mid C_t)$

Results

- ▶ No loss of optimality in using decision strategies $g_t^i: (L_t^i, \Pi_t) \mapsto U_t^i$.
- ▶ Dynamic program: $V(\pi) = \min_{\gamma} \mathbb{E}[R_t + \beta V(\Pi_{t+1}) \mid \Pi_t = \pi, \Gamma_t = \gamma]$.

-
- ▶ Witsenhausen, "Separation of estimation and control," Proc IEEE, 1971.
 - ▶ Nayyar, Mahajan and Teneketzis, "Optimal control strategies in delayed sharing information structures," IEEE TAC 2011.

Example 2: Control sharing information structure

Dynamics $X_{t+1}^i = f^i(X_t^i, \mathbf{u}_t, W_t^i)$, where $\mathbf{u}_t = (u_t^1, \dots, u_t^n)$.

Information **Original** : $I_t^i = \{X_{1:t}^i, \mathbf{u}_{1:t-1}\}$
structure **Using p-by-p approach:** $\tilde{I}_t^i = \{X_t^i, \mathbf{u}_{1:t-1}\}.$

Example 2: Control sharing information structure

Dynamics $X_{t+1}^i = f^i(X_t^i, \mathbf{U}_t, W_t^i), \quad \text{where } \mathbf{U}_t = (U_t^1, \dots, U_t^n).$

Information **Original** : $I_t^i = \{X_{1:t}^i, \mathbf{U}_{1:t-1}\}$
structure **Using p-by-p approach:** $\tilde{I}_t^i = \{X_t^i, \mathbf{U}_{1:t-1}\}.$

Common info.: $C_t = \mathbf{U}_{1:t-1}, \quad \text{Local Info.: } L_t^i = X_t^i, \quad \text{Prescriptions: } \Gamma_t^i: X_t^i \mapsto U_t^i$

Information Define $\Xi_t^i(x) = \mathbb{P}(X_t^i = x \mid \mathbf{U}_{1:t-1}).$

State Then $\Xi_t = (\Xi_t^1, \dots, \Xi_t^n)$ is an information state.

Results ▶ No loss of optimality in using decision strategies $g_t^i: (X_t^i, \Xi_t) \mapsto U_t^i.$

▶ Dynamic program: $V(\xi) = \min_{\gamma} \mathbb{E}[R_t + \beta V(\Xi_{t+1}) \mid \Xi_t = \xi, \Gamma_t = \gamma].$

Example 3: Mean-field sharing information structure

Dynamics $X_{t+1}^i = f_t(X_t^i, U_t^i, M_t, W_t^i),$ where $M_t = \sum_{i=1}^n \delta_{X_t^i}.$

Information structure $I_t^i = \{X_t^i, M_{1:t}\},$ and assume identical decision rules.

Example 3: Mean-field sharing information structure

Dynamics $X_{t+1}^i = f_t(X_t^i, U_t^i, M_t, W_t^i)$, where $M_t = \sum_{i=1}^n \delta_{X_t^i}$.

Information structure $I_t^i = \{X_t^i, M_{1:t}\}$, and assume identical decision rules.

Common info.: $C_t = M_{1:t}$, Local info.: $L_t^i = X_t^i$, Prescriptions: $\Gamma_t: X_t^i \mapsto U_t^i$.

Information state Due to the symmetry of the system, M_t is an information-state.

- Results**
- ▶ No loss of optimality in using decision strategies: $g_t^i(X_t^i, M_t)$.
 - ▶ Dynamic program: $V(m) = \min_{\gamma} \mathbb{E}[R_t + \beta V(M_{t+1}) \mid M_t = m, \Gamma_t = \gamma]$
 - ▶ Size of state space = $\text{poly}(n)$; Size of action space \mathcal{U}^x .

What if the shared information is empty?
The designer's approach

An example: Players with finite memory

Dynamics $X_{t+1} = f_t(X_t, U_t, W_t), \quad Y_t = h_t(X_t, N_t).$

Information structure $I_t = \{Y_t, M_t\}$ **Simplest non-classical information structure**
 $[U_t, M_{t+1}] = g_t(Y_t, M_t)$

-
- ▶ Witsenhausen, "A standard form for sequential stochastic control," Math. Sys. Theory, 1973.
 - ▶ Mahajan, "Sequential decomposition of sequential teams," PhD Thesis, 2008.

An example: Players with finite memory

Dynamics $X_{t+1} = f_t(X_t, U_t, W_t), \quad Y_t = h_t(X_t, N_t).$

Information structure $I_t = \{Y_t, M_t\}$ **Simplest non-classical information structure**
 $[U_t, M_{t+1}] = g_t(Y_t, M_t)$

Common info.: $C_t = \emptyset$, Local info.: $L_t = (Y_t, M_t)$, Prescriptions: $g_t: (Y_t, M_t) \mapsto U_t$.

Information state $\Pi_t = \mathbb{P}(X_t, M_t \mid g_{1:t-1})$

Results ▶ Dynamic program: $V(\pi) = \min_g \mathbb{E}[R_t + \beta V(\Pi_{t+1}) \mid \Pi_t = \pi, g_t = g]$

▶ Cannot show that time-homogeneous strategies are optimal!

▶ Witsenhausen, "A standard form for sequential stochastic control," Math. Sys. Theory, 1973.

▶ Mahajan, "Sequential decomposition of sequential teams," PhD Thesis, 2008.

Final Thoughts

Relation to game-theory

- ▶ Generalizes to Markov perfect equilibrium in stochastic games with asymmetric information (Nayyar, Gupta, Langbort, Başar, 2014).
- ▶ Implications for dynamic mechanism design.

Is common information (or PHS) a realistic assumption?

- ▶ Arises naturally in certain applications.
- ▶ Use (a faster time-scale) consensus dynamics to generate common information
- ▶ Provide upper and lower bounds

Are there good numerical algorithms?

- ▶ Are there POMDP algorithms for large action spaces?
- ▶ Is there some structure in the DP that can be exploited?

Interesting variations

- ▶ ε common-information
- ▶ Approximation techniques
- ▶ Reinforcement learning
- ▶ Other information structures (sparse structures)?

References

Nayyar, Mahajan and Teneketzis, “Decentralized stochastic control with partial history sharing: A common information approach,” IEEE TAC 2013.

Mahajan, Nayyar, and Teneketzis, “Identifying tractable decentralized problems on the basis of information structures”, Allerton 2008.

Nayyar, “Sequential Decision-Making in Decentralized systems,” PhD Thesis, Univ of Michigan, 2011.

Nayyar, Mahajan and Teneketzis, “Optimal control strategies in delayed sharing information structures,” IEEE TAC 2011.

Mahajan, “Optimal decentralized control of coupled subsystems with control sharing,” IEEE TAC 2013.

Arabneydi and Mahajan, “Team optimal control of coupled subsystems with mean field sharing,” CDC 2014.

Mahajan and Mannan, “Decentralized Stochastic Control,” Annals of OR, (in print).