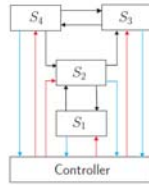


Control of Networked Markov Decision Systems with Delays

Sachin Adlakha, Sanjay Lall and Andrea Goldsmith

Introduction

- Feedback control of networked systems.
- Controller has *delayed* access to system states.
- Each subsystem a *Markov Decision Process*.



We show

An optimal controller exists with *finite memory*

- We call the controller *banded* by analogy with matrix case.
- Resulting dynamic program is computationally tractable.

Main Result

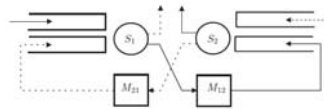
The optimal controller is of the form $u_t = \mu_t^{\text{opt}}(i_t^{\text{mem}})$ where

$$i_t^{\text{mem}} = \left(u_{t-d_i:t-1}^i, x_{t-N_i-b_i:t-N_i}^i \mid i = 1, 2, \dots, n \right)$$

- Optimal controller only needs *part* of the observation history
- Required history is i_t^{mem}
- Required history is *finite*, length does not depend on t
- Constants b_k, d_k , called *bandwidths* depend on graph structure and delays, not on dynamics

Example : Interconnected Queues

$$\begin{aligned} x^1(t+1) &= f^1(x^1(t), x^2(t-M_{21}), u^1(t), w^1(t)), \\ x^2(t+1) &= f^2(x^2(t), x^1(t-M_{12}), u^2(t), w^2(t)) \end{aligned}$$

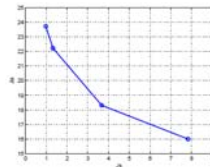


- The state cost is

$$g_s(x^1(t), x^2(t)) = (x_R^1(t) + x_B^1(t) + x_R^2(t) + x_B^2(t))^2.$$

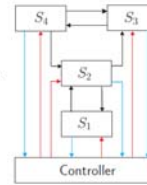
- The action cost is

$$g_a(u^1(t), u^2(t)) = (u^1(t) + 1 + u^2(t) + 1)^2.$$



Motivation

- Networked Markov decision processes are used to model
 - Distributed vehicle coordination.
 - Scheduling over multiple servers.
 - Network of interacting queues.



- Delays cause the system to be *partially observable*.

Provide sufficient information state for such systems.

Main Result

The optimal controller is of the form $u_t = \mu_t^{\text{opt}}(i_t^{\text{mem}})$ where

$$i_t^{\text{mem}} = \left(u_{t-d_i:t-1}^i, x_{t-N_i-b_i:t-N_i}^i \mid i = 1, 2, \dots, n \right)$$

The bandwidths b_k and d_k are

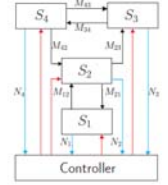
$$b_k = \max \{ d_k, d_s + M_{ks} \mid s \in \mathcal{O}^k \} - N_k$$

$$d_k = \max \{ N_k, N_s - 1 - M_{sk} \mid s \in \mathcal{I}^k \}$$

- \mathcal{I}^k is the set of vertices with incoming edges into vertex k .
- \mathcal{O}^k is the set of vertices with outgoing edges from vertex k .

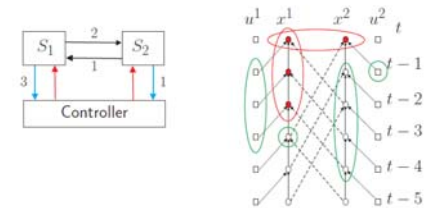
Model

$$\begin{aligned} x_{t+1}^1 &= f^1(x_t^1, u_t^1, w_t^1, x_{t-M_{21}}^2) \\ x_{t+1}^2 &= f^2(x_t^2, u_t^2, w_t^2, x_{t-M_{12}}^1, x_{t-M_{42}}^4) \\ x_{t+1}^3 &= f^3(x_t^3, u_t^3, w_t^3, x_{t-M_{23}}^2, x_{t-M_{43}}^4) \\ x_{t+1}^4 &= f^4(x_t^4, u_t^4, w_t^4, x_{t-M_{34}}^3) \end{aligned}$$



- x_t^i is the state of subsystem S_i at time t
- M_{ij} is the delay from subsystem S_i to subsystem S_j .
- N_i is the delay in receiving observations from subsystem S_i .
- The controller gets states from every subsystem.

Connections to Bayesian Networks



It is easy to check that

i_t^{mem} forms the Markov blanket.

Conclusions

Networked MDPs

- Optimal control synthesis
- Centralized control, distributed delayed state feedback
- Finite state systems

An optimal controller exists with finite memory.

- Bandwidths depend only on the graph structure and the delays.
- The bands forms a Markov blanket for unknown states.

- Minimize the infinite horizon discounted cost

$$\begin{aligned} J &= \mathbb{E} \left(\sum_{t=0}^{\infty} \beta^t ((1-\alpha)g_s + \alpha g_a) \right), \\ &= (1-\alpha)J_s + \alpha J_a \end{aligned}$$

- The propagation and measurement delays are

$$M_{12} = 2, \quad M_{21} = 1, \quad \text{and} \quad N_1 = 2, \quad N_2 = 1$$

The bands for the optimal controller are

$$b_1 = 1, \quad b_2 = 2 \quad \text{and} \quad d_1 = 2, \quad d_2 = 1.$$