

Teaching the correct policy to ENMPC

Adaptive Economic NMPC and Reinforcement Learning

Sébastien Gros

Cybernetic, NTNU
Electrical Engineering, Chalmers

TUM

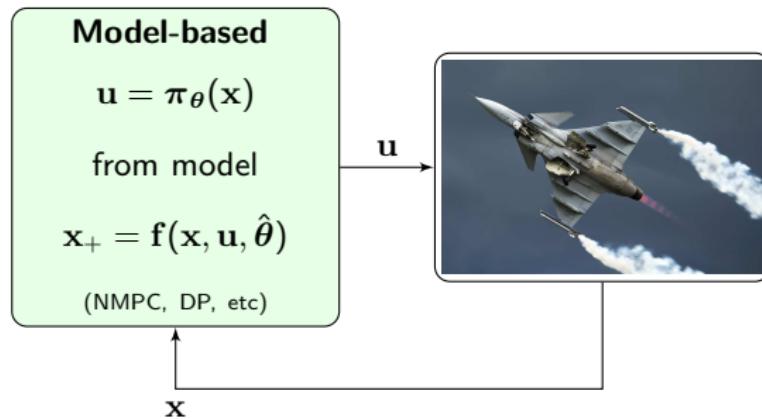
Outline

- 1 Introduction
- 2 Modification of the ENMPC scheme
- 3 Learning the Optimal ENMPC
- 4 RL for NMPC in practice
- 5 Simple example
- 6 Conclusions

Outline

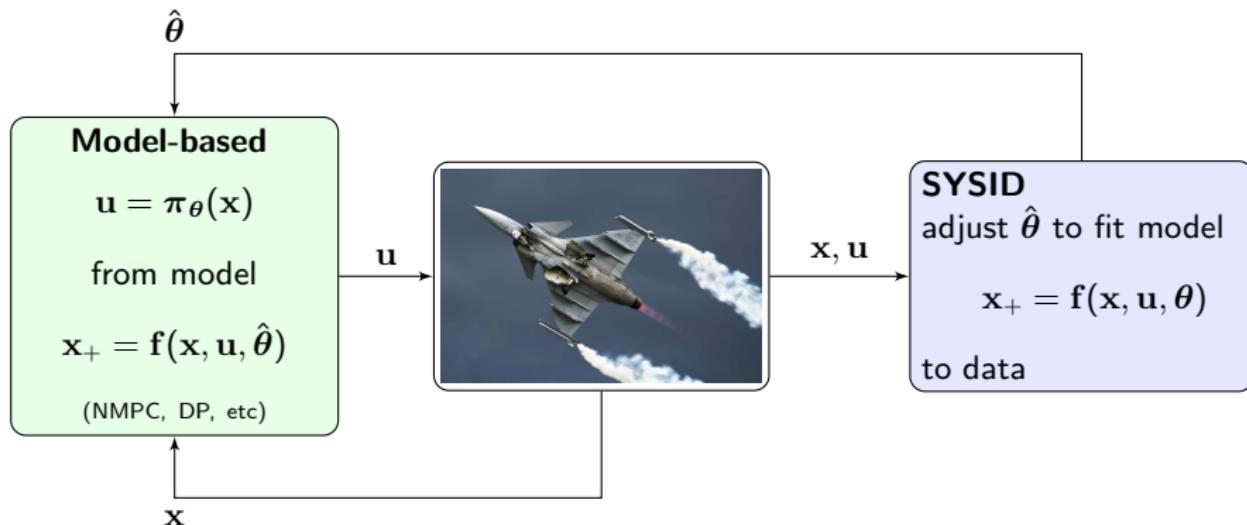
- 1 Introduction
- 2 Modification of the ENMPC scheme
- 3 Learning the Optimal ENMPC
- 4 RL for NMPC in practice
- 5 Simple example
- 6 Conclusions

SYSID & Model-based optimal control



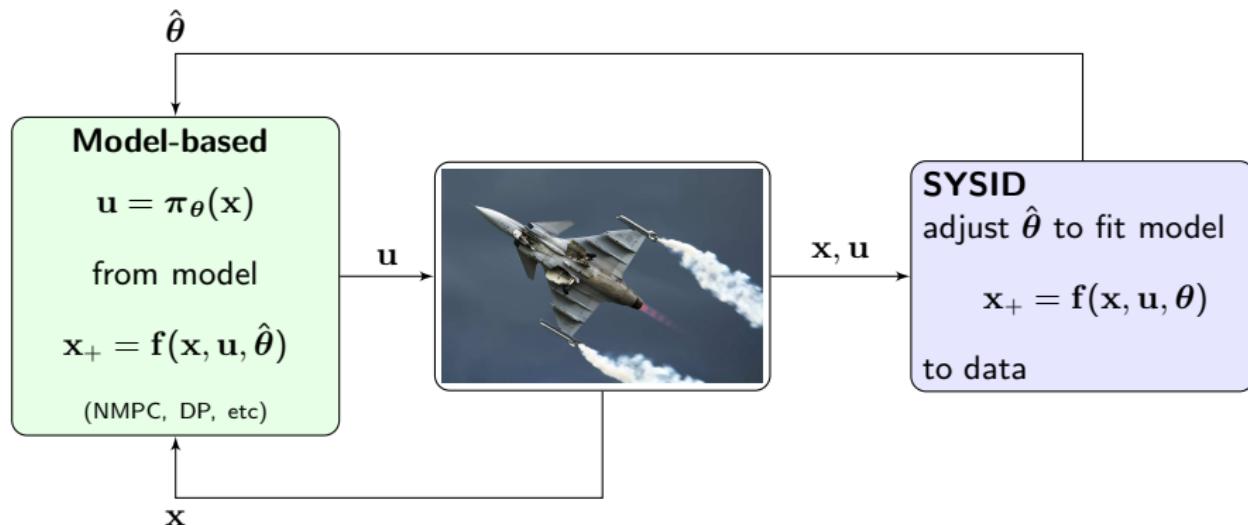
Assume we have a given performance index (i.e. a stage cost L to minimized)

SYSID & Model-based optimal control



Assume we have a given performance index (i.e. a stage cost L to minimized)

SYSID & Model-based optimal control



Assume we have a given performance index (i.e. a stage cost L to minimized)

Does this work?

- Not necessarily... problem: model may not able to capture the real system
- E.g. what is the best way of representing stochasticity using f ?
- Can degrade performance compared to keeping initial $\hat{\theta}$

Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .

Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

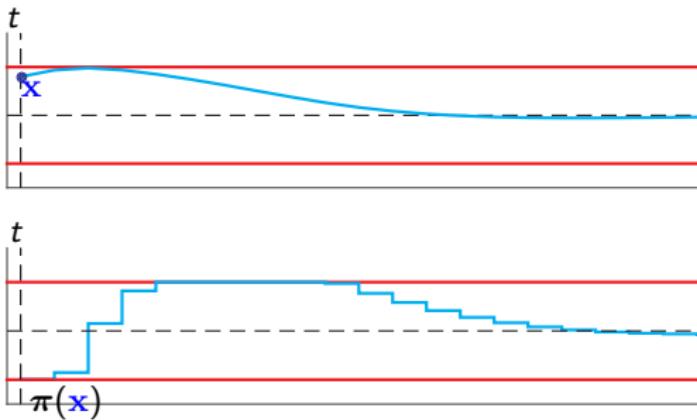
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

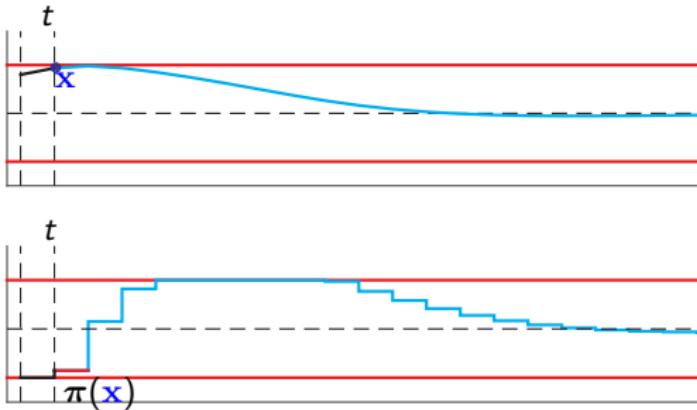
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

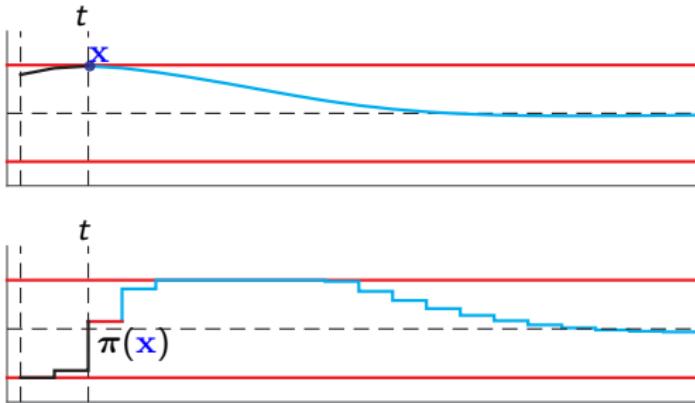
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

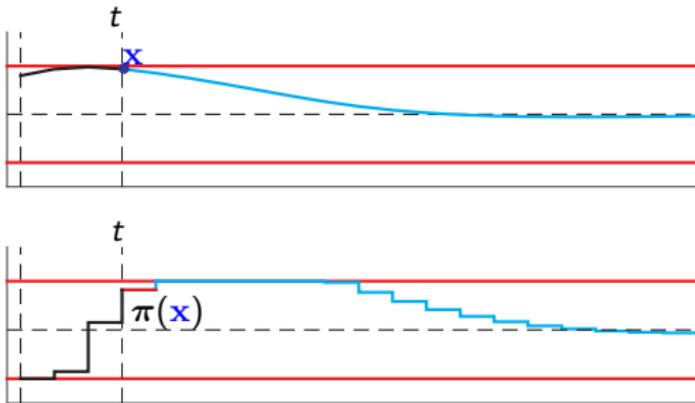
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

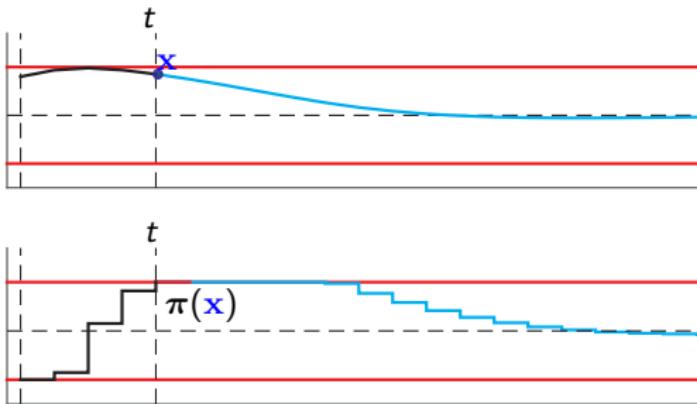
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

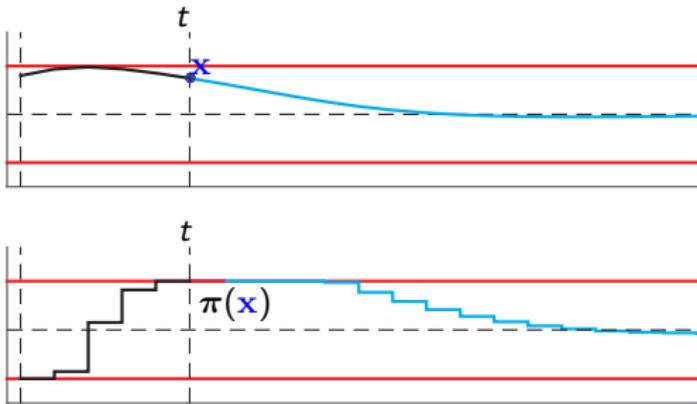
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

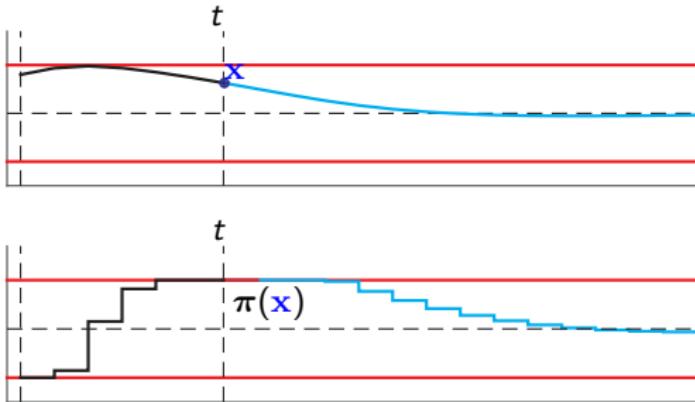
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

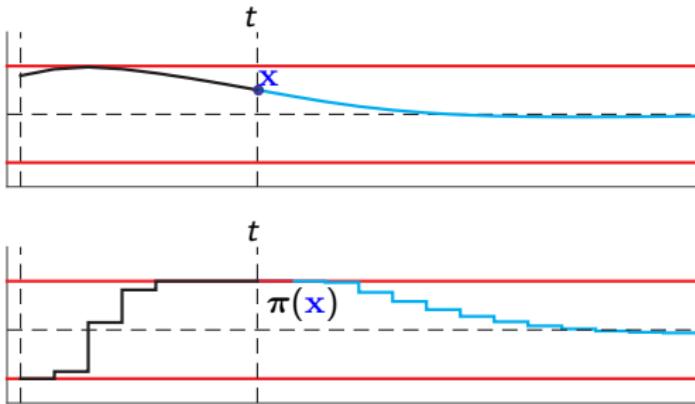
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

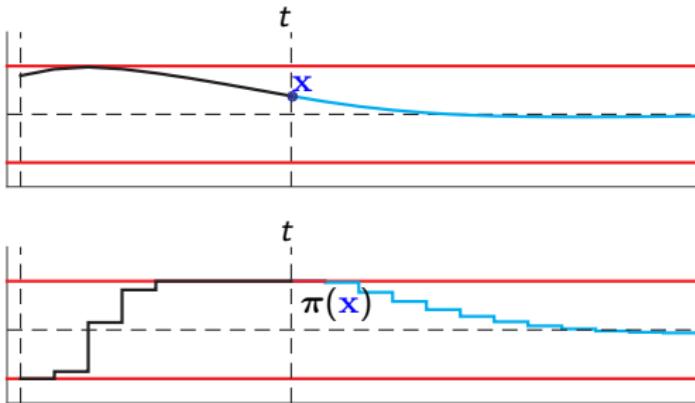
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

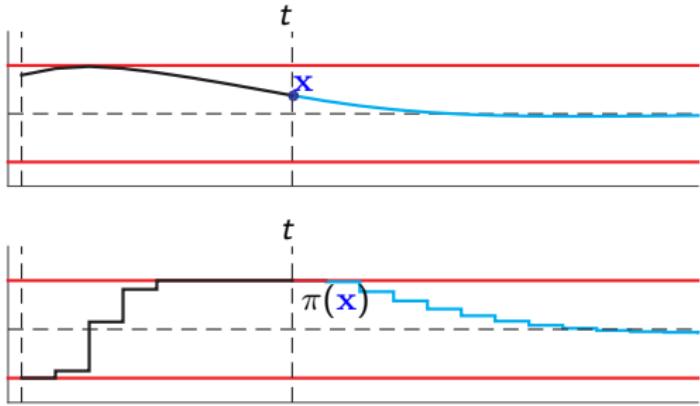
$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{u}} \quad & \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right] \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0 \\ & \mathbf{x}_0 = \mathbf{x} \end{aligned}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

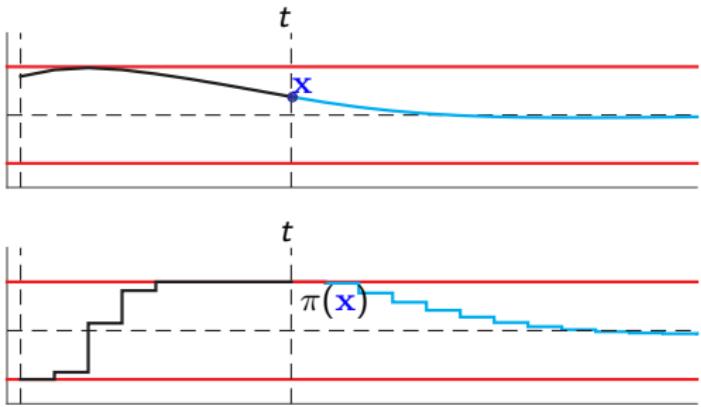
Always update \mathbf{u} using the latest information \mathbf{x}

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned} \mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \quad & \gamma^N T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t.} \quad & \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) \\ & \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0 \\ & \mathbf{x}_0 = \mathbf{x} \end{aligned}$$

for given model \mathbf{f} .



Optimal control

Minimize “long-term” cost

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{u}} \quad & \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right] \\ \text{s.t.} \quad & \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0 \\ & \mathbf{x}_0 = \mathbf{x} \end{aligned}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned} \mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \quad & \gamma^N T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t.} \quad & \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) \\ & \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0 \\ & \mathbf{x}_0 = \mathbf{x} \end{aligned}$$

for given model \mathbf{f} .

Remarks

- usually $\gamma = 1$
- Generic $L \rightarrow$ Economic NMPC

Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned} \mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } &\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) \\ &\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0 \\ &\mathbf{x}_0 = \mathbf{x} \end{aligned}$$

for given model \mathbf{f} .

Remarks

- usually $\gamma = 1$
- Generic $L \rightarrow$ Economic NMPC

User manual

- Fit model \mathbf{f} to real dynamics (SYSID)
- Good choice of T

and hope that $\pi \approx \pi_*$

Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Model-based optimal control - Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned} \mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } &\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) \\ &\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0 \\ &\mathbf{x}_0 = \mathbf{x} \end{aligned}$$

for given model \mathbf{f} .

Problem:

- model \mathbf{f} may not be able to capture real dynamics
- real dynamics are often stochastic

nothing guarantees that π from NMPC yields good performances

Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

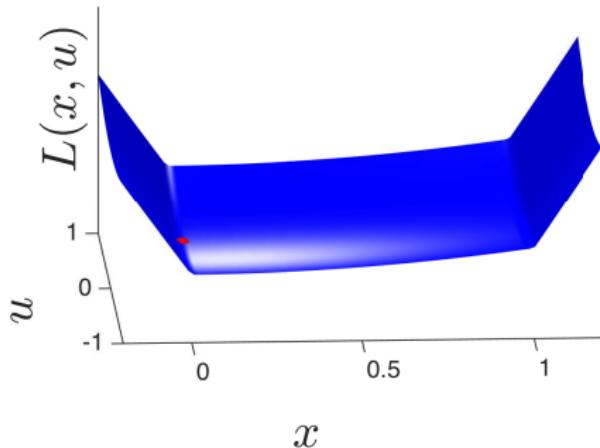
Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Our example

$$x_+ = x + 0.1u + e, \quad e \sim \mathcal{U}([-0.1, 0])$$



- Quadratic, minimum is at $x = 0, u = 0$
- Strong penalty for $x \notin [0, 1]$ and $u \in [-1, 1]$
- Noise e “pushing” x “to the left”

SYSID with deterministic model:

$$\hat{x}_+ = x + 0.1u - 0.05$$

Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

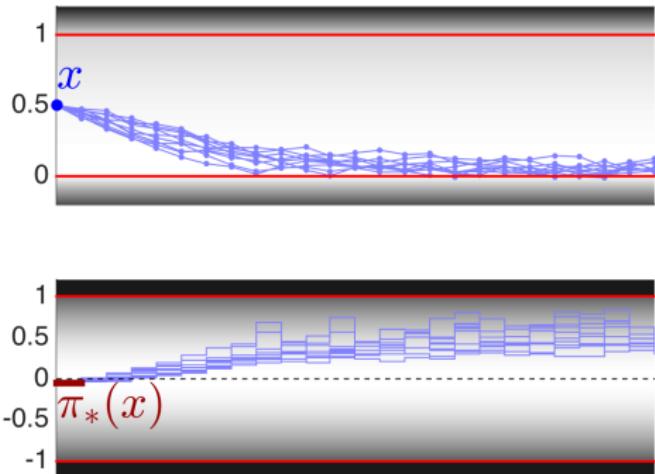
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

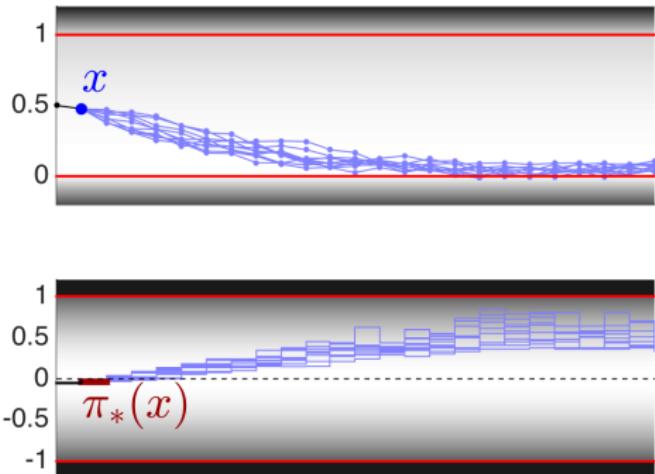
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

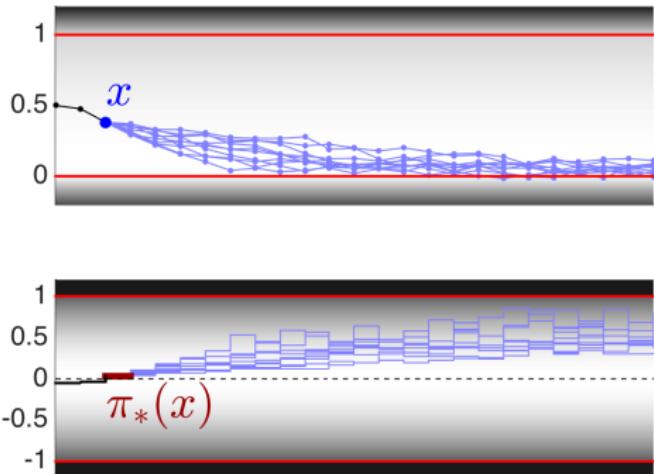
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

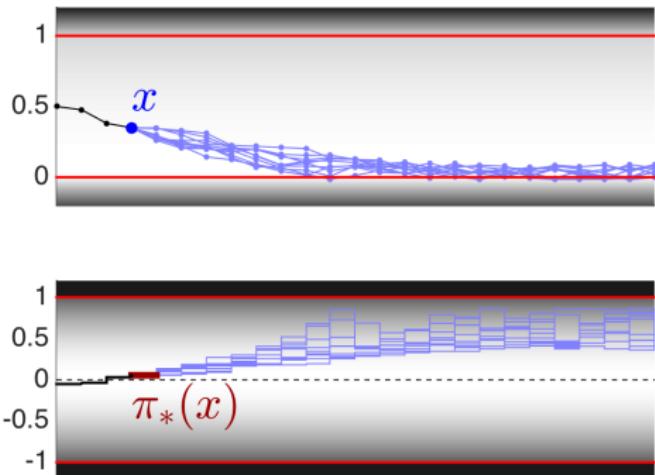
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

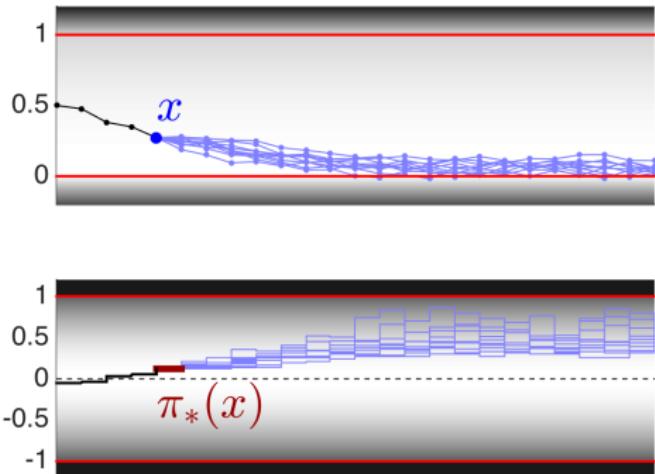
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

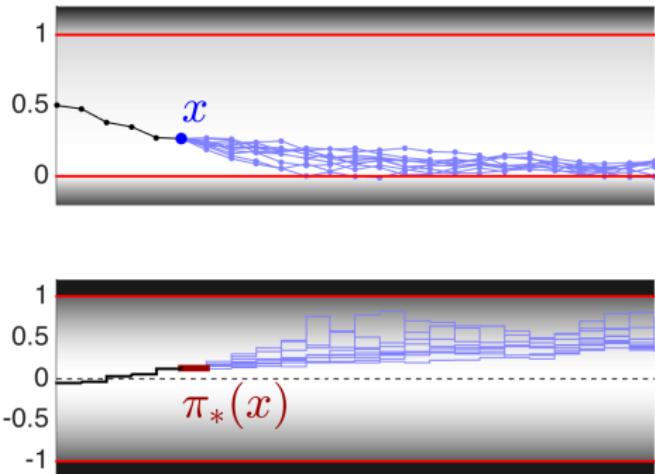
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

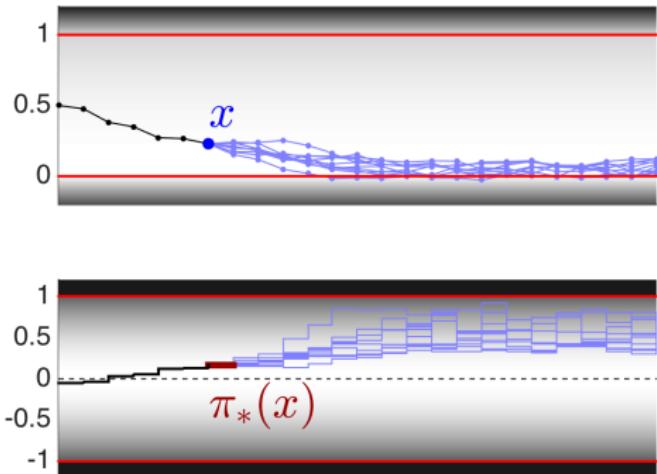
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

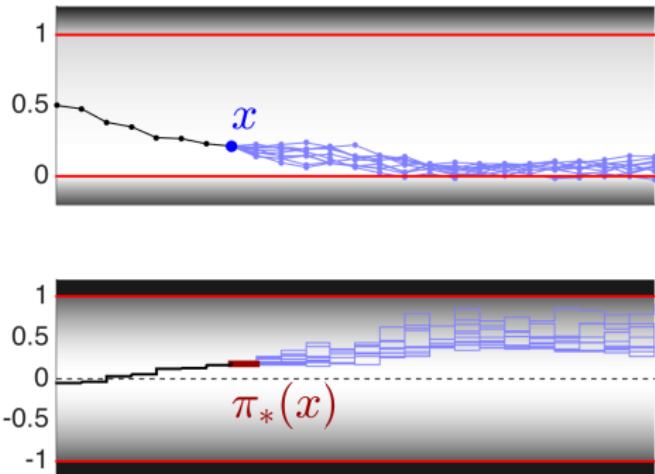
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

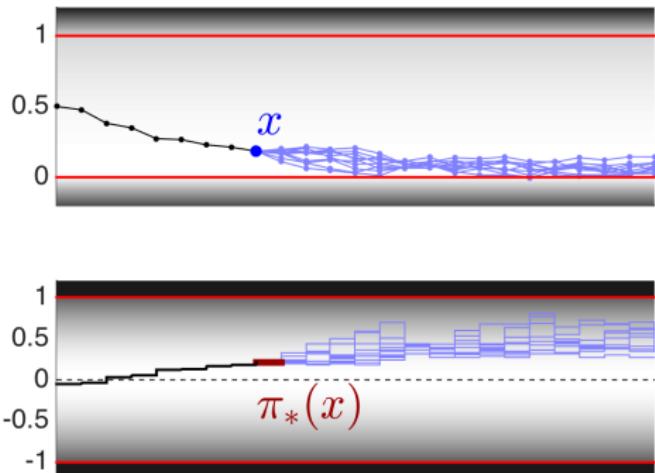
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

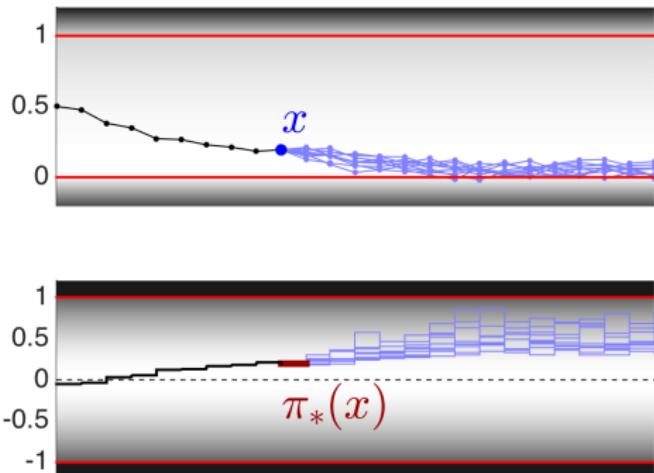
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



Optimal control

Minimize “long-term” cost

$$\min_{\mathbf{x}, \mathbf{u}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \right]$$

$$\text{s.t. } \mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

given initial condition \mathbf{x} .

Under real dynamics:

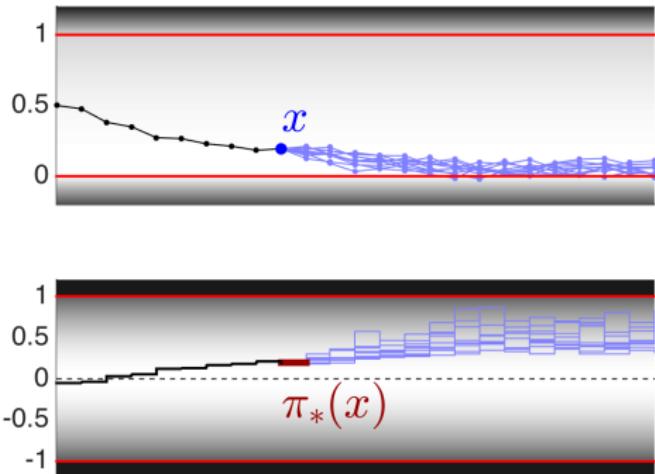
$$\mathbb{P}[\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k]$$

Yields optimal **policy**:

$$\mathbf{u} = \pi_*(\mathbf{x})$$

Always update \mathbf{u} using the latest information \mathbf{x}

Optimal trajectories (policy π_* (\mathbf{x}))
(dark = no good)



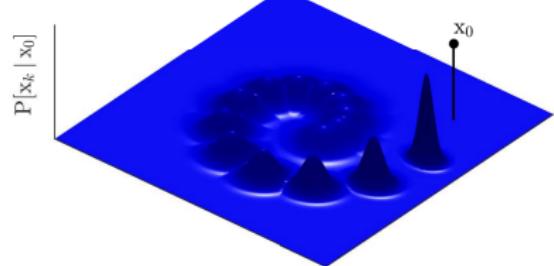
This is a Markov Decision Process (MDP):

- Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$
- Stage cost
- Discount $\gamma \in]0, 1]$

A generic framework - Markov Decision Processes

Markov Decision Process (MDP):

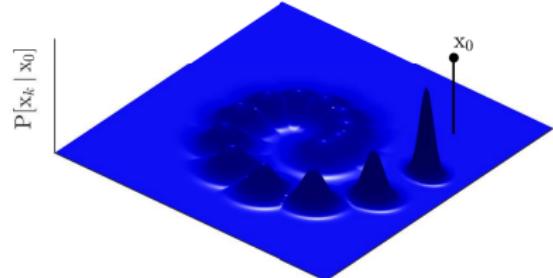
- Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$
- Stage cost: $\ell(\mathbf{x}, \mathbf{u}) \in \mathbb{R}$
- Discount factor: $\gamma \in [0, 1]$



A generic framework - Markov Decision Processes

Markov Decision Process (MDP):

- Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$
- Stage cost: $\ell(\mathbf{x}, \mathbf{u}) \in \mathbb{R}$
- Discount factor: $\gamma \in [0, 1]$



Constraints: $\mathbf{h}(\mathbf{x}, \mathbf{u}) \leq 0$, then define

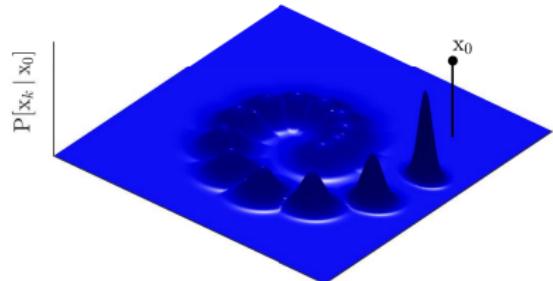
$$\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } \mathbf{h}(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

where $L(\mathbf{x}, \mathbf{u})$ is the finite stage cost

A generic framework - Markov Decision Processes

Markov Decision Process (MDP):

- Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$
- Stage cost: $\ell(\mathbf{x}, \mathbf{u}) \in \mathbb{R}$
- Discount factor: $\gamma \in [0, 1]$



Constraints: $h(\mathbf{x}, \mathbf{u}) \leq 0$, then define

$$\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

where $L(\mathbf{x}, \mathbf{u})$ is the finite stage cost

Optimal policy $\mathbf{u} = \pi_*(\mathbf{x})$ yields value function:

$$V_*(\mathbf{x}_0) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \ell(\mathbf{x}_k, \pi_*(\mathbf{x}_k)) \mid \mathbf{x}_0, \pi_* \right]$$

(can take ∞ values)

Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Solving MDPs - Dynamic Programming & the Bellman equations

Real system

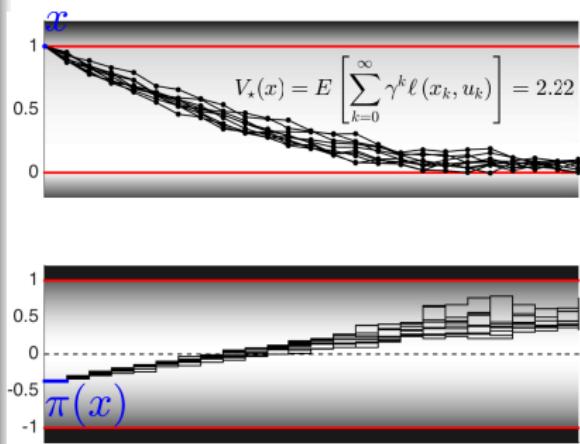
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

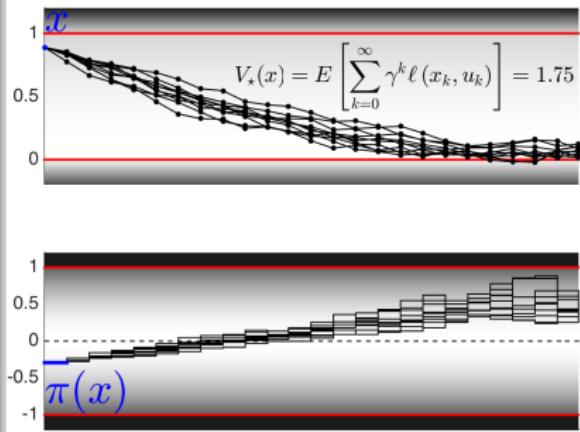
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

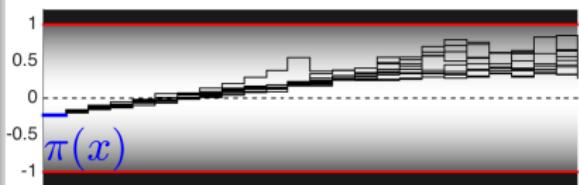
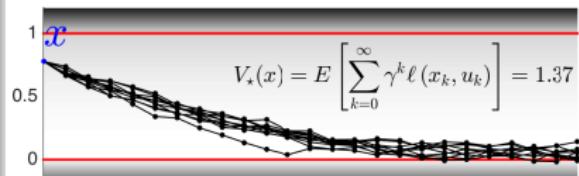
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

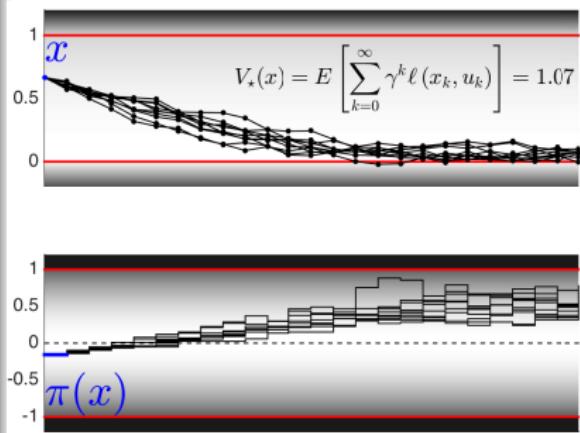
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

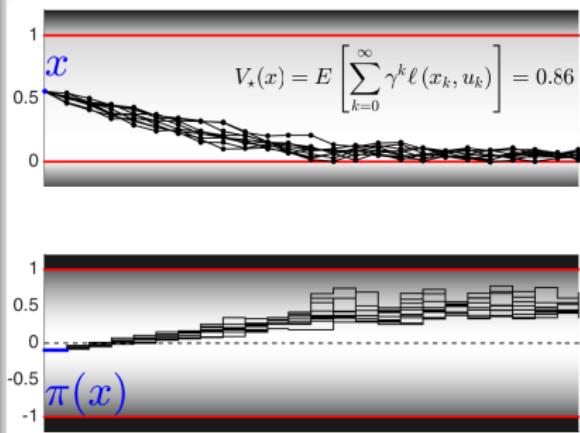
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

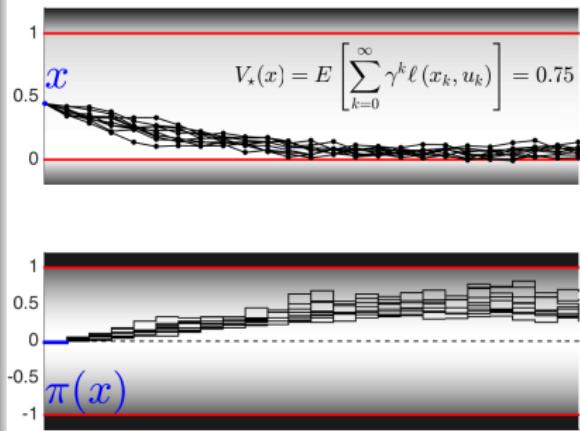
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

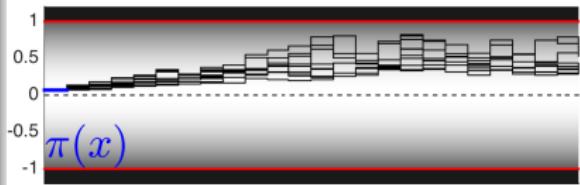
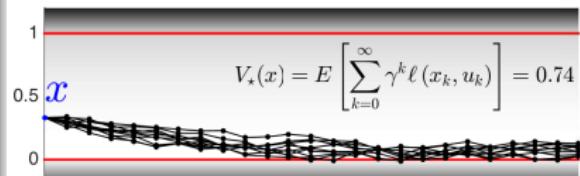
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

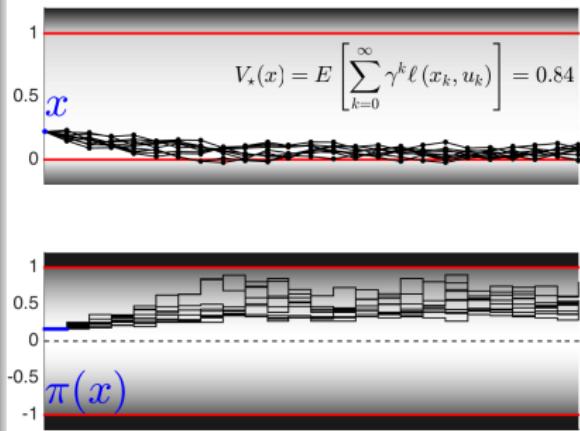
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

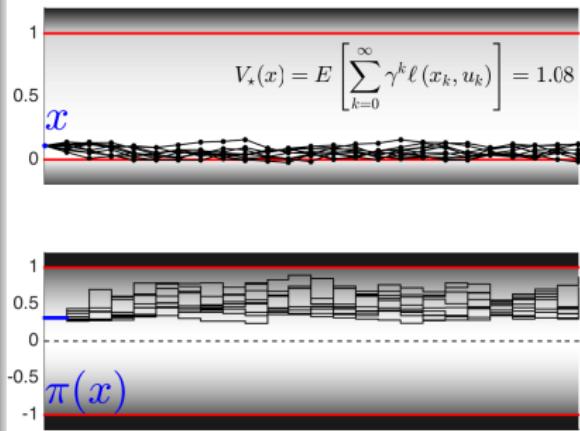
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

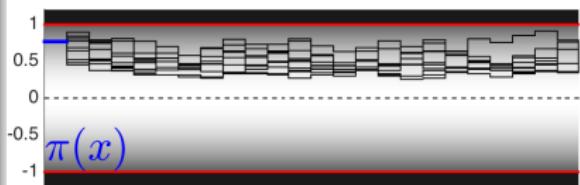
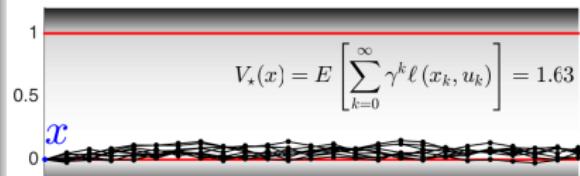
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

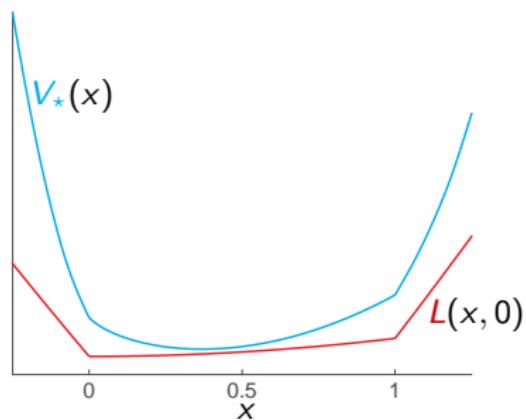
Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

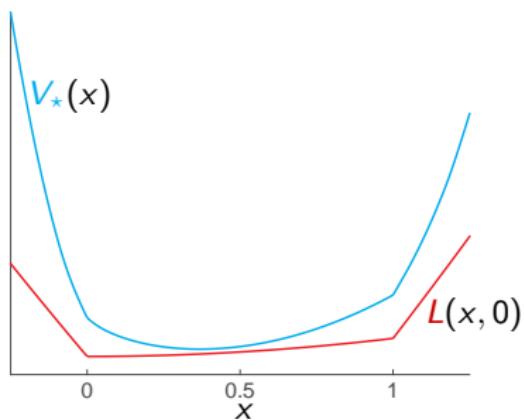
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

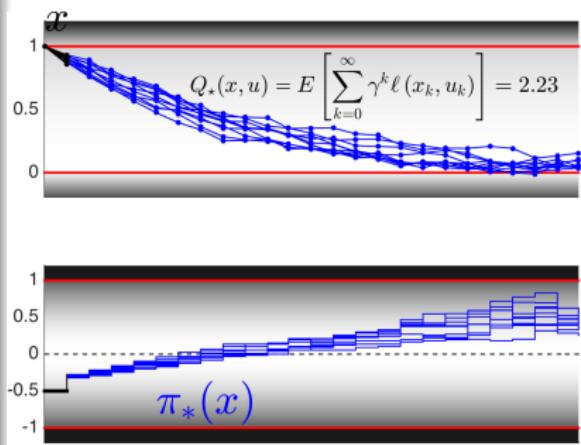
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

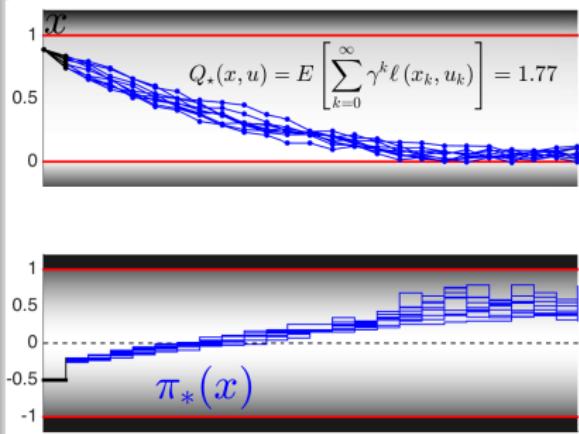
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

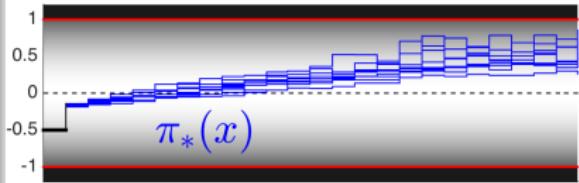
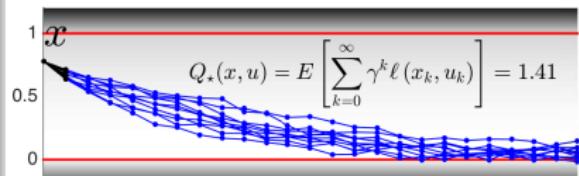
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

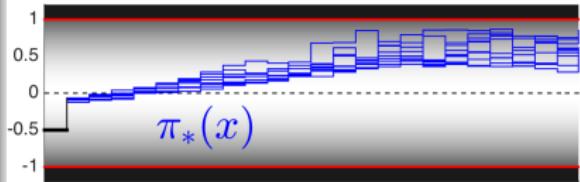
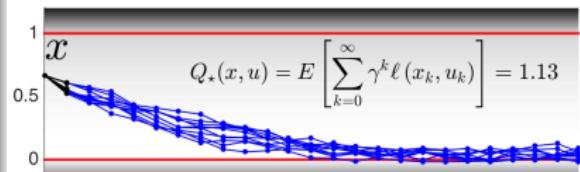
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

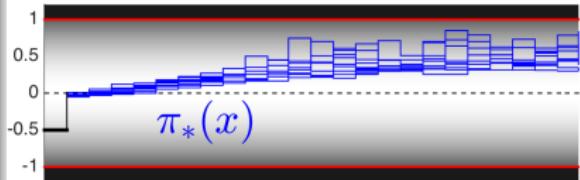
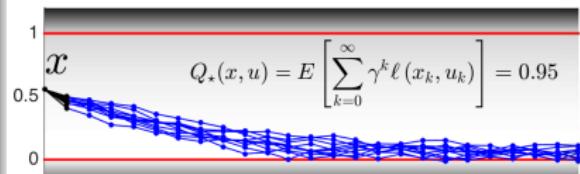
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

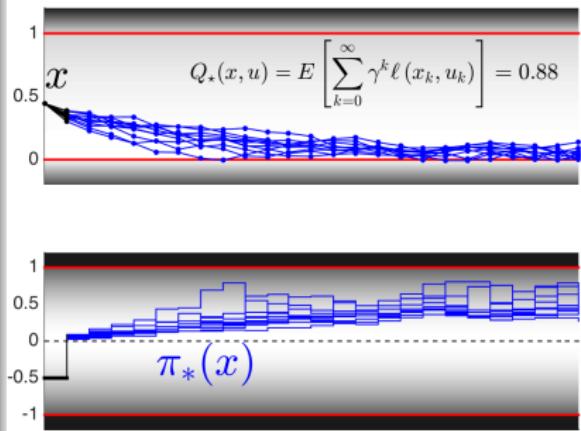
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

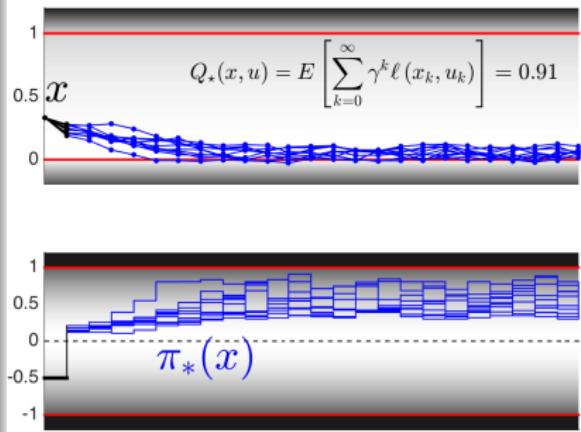
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

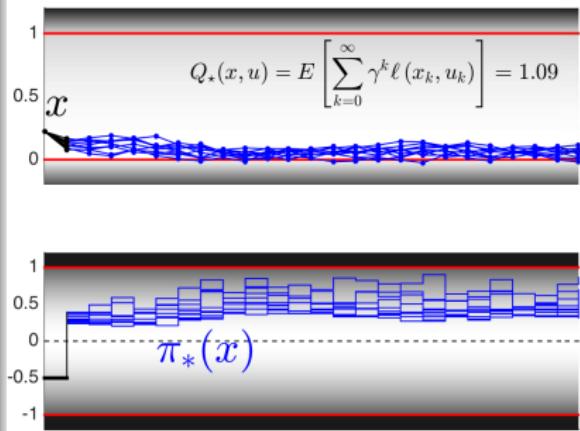
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

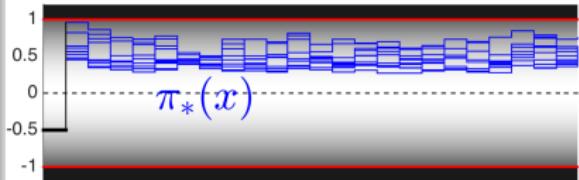
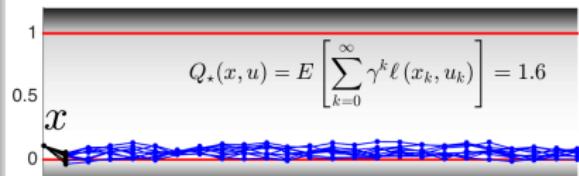
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

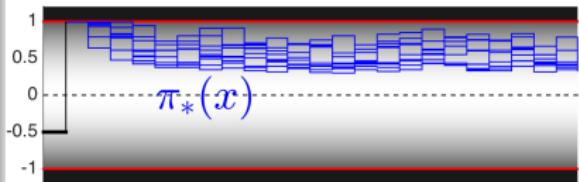
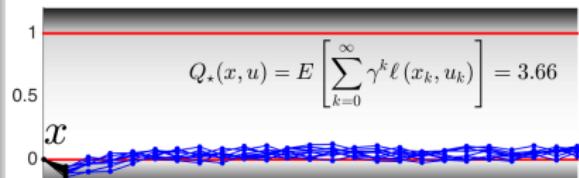
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

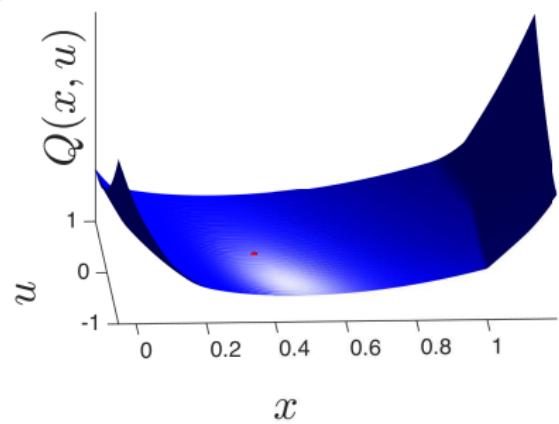
Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

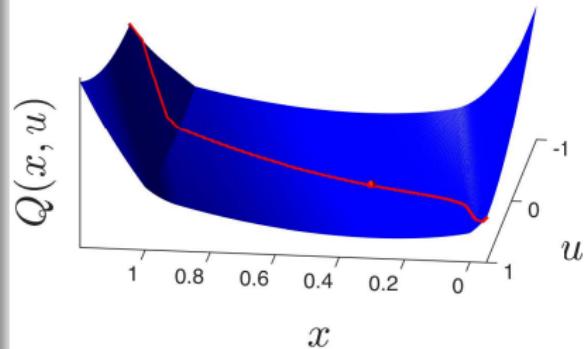
$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Relationship $V_* \leftrightarrow Q_*$:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}) = Q_*(\mathbf{x}, \pi_*(\mathbf{x}))$$

Our trivial example



Solving MDPs - Dynamic Programming & the Bellman equations

Real system

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Bellman optimality equations:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = \ell(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

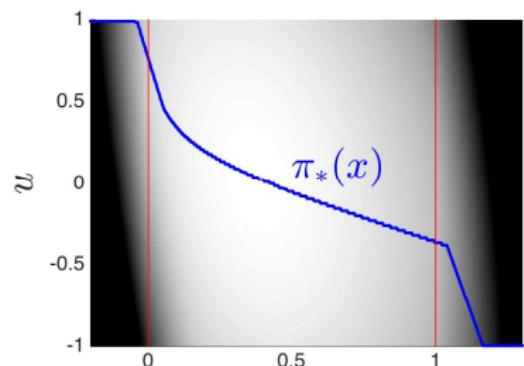
Relationship $V_* \leftrightarrow Q_*$:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}) = Q_*(\mathbf{x}, \pi_*(\mathbf{x}))$$

Optimal policy:

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Our trivial example



Outline

- 1 Introduction
- 2 Modification of the ENMPC scheme
- 3 Learning the Optimal ENMPC
- 4 RL for NMPC in practice
- 5 Simple example
- 6 Conclusions

Modification of the ENMPC scheme

Real system - MDP

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Value functions: $Q_*(\mathbf{x}, \mathbf{u}), V_*(\mathbf{x})$

Optimal policy $\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$h(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

Modification of the ENMPC scheme

Real system - MDP

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Value functions: $Q_*(\mathbf{x}, \mathbf{u}), V_*(\mathbf{x})$

Optimal policy $\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

Modification of the ENMPC scheme

Real system - MDP

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Value functions: $Q_*(\mathbf{x}, \mathbf{u}), V_*(\mathbf{x})$

Optimal policy $\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

Can we gain from modifying the cost and constraints?†

† idea coming from Real-Time Optimization

Modification of the ENMPC scheme

Real system - MDP

Dynamics: $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Stage cost: $\ell(\mathbf{x}, \mathbf{u}) = \begin{cases} L(\mathbf{x}, \mathbf{u}) & \text{if } h(\mathbf{x}, \mathbf{u}) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Value functions: $Q_*(\mathbf{x}, \mathbf{u}), V_*(\mathbf{x})$

Optimal policy $\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{L}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

Can we gain from modifying the cost and constraints?†

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong?

† idea coming from Real-Time Optimization

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \quad \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Theorem

Under some assumptions
there is a $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ such that

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

$$V(\mathbf{x}) = V_*(\mathbf{x})$$

$$Q(\mathbf{x}, \mathbf{u}) = Q_*(\mathbf{x}, \mathbf{u})$$

π, V, Q delivered by ENMPC

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Theorem

Under some assumptions
there is a $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ such that

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

$$V(\mathbf{x}) = V_*(\mathbf{x})$$

$$Q(\mathbf{x}, \mathbf{u}) = Q_*(\mathbf{x}, \mathbf{u})$$

π, V, Q delivered by ENMPC

Bottom line: modifying the cost and constraints of the (E)NMPC may help gaining performance when the model is wrong.

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \quad \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Theorem

There is a $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi = \pi_*, Q = Q_*, V = V_*$$

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Theorem

There is a $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi = \pi_*, Q = Q_*, V = V_*$$

on $\mathbf{x} \in \mathcal{S}$ with:

$$\mathcal{S} = \{\mathbf{x}_0 \mid V_*(\mathbf{x}_k) \text{ finite } \forall k\}$$

where $\mathbf{x}_{0,\dots,\infty}$ is given by

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \pi_*(\mathbf{x}_k))$$

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{T}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{L}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{T}, \hat{L}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Theorem

There is a $\hat{T}, \hat{L}, \hat{\mathbf{h}}$ s.t.

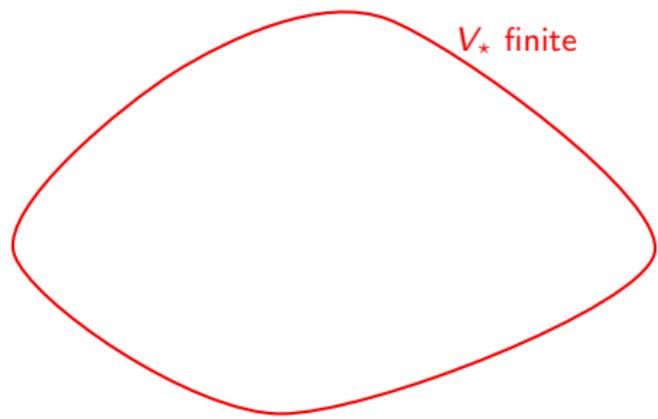
$$\pi = \pi_*, Q = Q_*, V = V_*$$

on $\mathbf{x} \in \mathcal{S}$ with:

$$\mathcal{S} = \{\mathbf{x}_0 \mid V_*(\mathbf{x}_k) \text{ finite } \forall k\}$$

where $\mathbf{x}_{0, \dots, \infty}$ is given by

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \pi_*(\mathbf{x}_k))$$



Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{T}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{L}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{T}, \hat{L}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Theorem

There is a $\hat{T}, \hat{L}, \hat{\mathbf{h}}$ s.t.

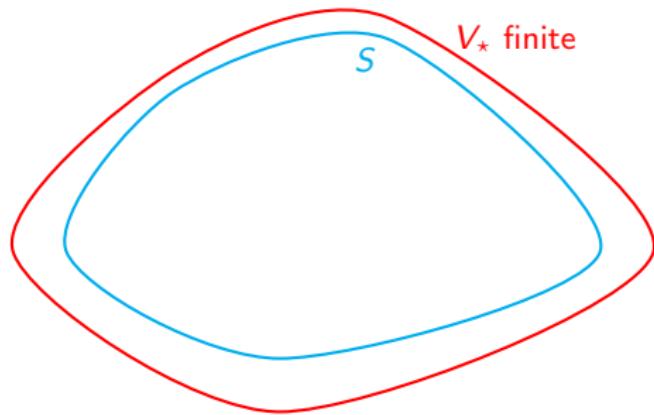
$$\pi = \pi_*, Q = Q_*, V = V_*$$

on $\mathbf{x} \in S$ with:

$$S = \{\mathbf{x}_0 \mid V_*(\mathbf{x}_k) \text{ finite } \forall k\}$$

where $\mathbf{x}_{0, \dots, \infty}$ is given by

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \pi_*(\mathbf{x}_k))$$



Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Theorem

There is a $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

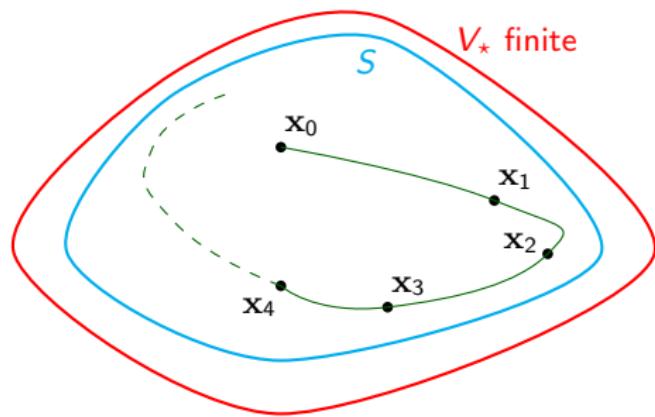
$$\pi = \pi_*, Q = Q_*, V = V_*$$

on $\mathbf{x} \in S$ with:

$$S = \{\mathbf{x}_0 \mid V_*(\mathbf{x}_k) \text{ finite } \forall k\}$$

where $\mathbf{x}_{0,\dots,\infty}$ is given by

$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \pi_*(\mathbf{x}_k))$$



Forward-invariant set of the model dynamics under π_* within set where V_* is finite

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(x) = u_0^*$ where

$$\begin{aligned} \mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0 \end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(x) = \pi_*(x)$$

even though the model f is wrong ?

Theorem

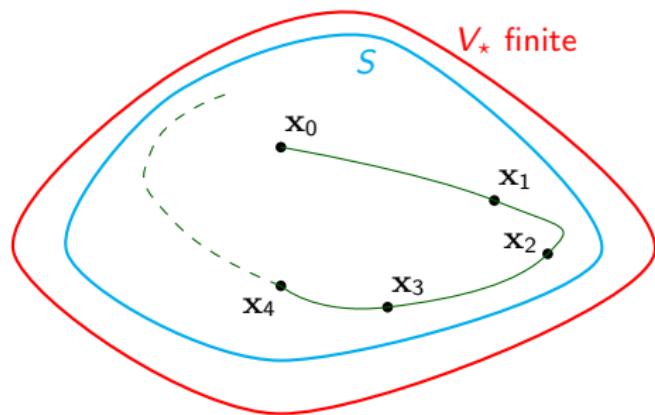
If $\hat{\mathcal{T}} = 0$, and if for k large

$$\gamma^k V_*(\mathbf{x}_k) \rightarrow 0$$

then there is $\hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(x) \rightarrow \pi_*(x) \text{ on } x \in S$$

as N chosen large



Forward-invariant set of the model dynamics under π_* within set where V_* is finite

Modification of the ENMPC scheme

Economic NMPC

Policy $\pi(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\begin{aligned}\mathbf{u}^*, \mathbf{x}^* &= \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \hat{\mathcal{T}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \hat{\mathcal{L}}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x} \\ \hat{\mathbf{h}}(\mathbf{x}_k, \mathbf{u}_k) &\leq 0\end{aligned}$$

Is there $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) = \pi_*(\mathbf{x})$$

even though the model \mathbf{f} is wrong ?

Theorem

If $\hat{\mathcal{T}} = 0$, and if for k large

$$\gamma^k V_*(\mathbf{x}_k) \rightarrow 0$$

then there is $\hat{\mathcal{L}}, \hat{\mathbf{h}}$ s.t.

$$\pi(\mathbf{x}) \rightarrow \pi_*(\mathbf{x}) \text{ on } \mathbf{x} \in \mathcal{S}$$

as N chosen large

Results also hold for

- (E)NMPC based on stochastic models (e.g. scenario tree, tube-based etc.)
- $\gamma = 1$, i.e. "classic" (E)NMPC formulations
- Deterministic real system

No free lunch: $\hat{\mathcal{T}}, \hat{\mathcal{L}}, \hat{\mathbf{h}}$ can be very complex

Outline

- 1 Introduction
- 2 Modification of the ENMPC scheme
- 3 Learning the Optimal ENMPC
- 4 RL for NMPC in practice
- 5 Simple example
- 6 Conclusions

Parametrization of the (E)NMPC scheme

NMPC policy: $\pi_{\theta}(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \mathcal{T}_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \mathcal{L}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

i.e. approximate $\hat{\mathcal{T}}$, $\hat{\mathcal{L}}$, $\hat{\mathbf{h}}$ and
give freedom in the model

Parametrization of the (E)NMPC scheme

NMPC policy: $\pi_{\theta}(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \mathcal{T}_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \mathcal{L}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

i.e. approximate $\hat{\mathcal{T}}$, $\hat{\mathcal{L}}$, $\hat{\mathbf{h}}$ and
give freedom in the model

- \mathcal{T}_{θ} , \mathcal{L}_{θ} , \mathbf{h}_{θ} are in theory enough
- Adjusting model \mathbf{f}_{θ} can only help “closing the gap” for limited parametrizations

Parametrization of the (E)NMPC scheme

NMPC policy: $\pi_\theta(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \mathcal{T}_\theta(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \mathcal{L}_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_\theta(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

i.e. approximate $\hat{\mathcal{T}}$, $\hat{\mathcal{L}}$, $\hat{\mathbf{h}}$ and give freedom in the model

- $\mathcal{T}_\theta, \mathcal{L}_\theta, \mathbf{h}_\theta$ are in theory enough
- Adjusting model \mathbf{f}_θ can only help “closing the gap” for limited parametrizations

(E)NMPC parameters θ ought to be selected to minimize

$$J(\theta) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \ell(\mathbf{x}_k, \pi_\theta(\mathbf{x}_k)) \right] \quad \text{where } \mathbf{x}_{1,\dots,\infty} \text{ result from } \pi_\theta \rightarrow \text{real system}$$

rather than via model fitting!!

Parametrization of the (E)NMPC scheme

NMPC policy: $\pi_\theta(\mathbf{x}) = \mathbf{u}_0^*$ where

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N \mathcal{T}_\theta(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \mathcal{L}_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_\theta(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

i.e. approximate $\hat{\mathcal{T}}$, $\hat{\mathcal{L}}$, $\hat{\mathbf{h}}$ and give freedom in the model

- $\mathcal{T}_\theta, \mathcal{L}_\theta, \mathbf{h}_\theta$ are in theory enough
- Adjusting model \mathbf{f}_θ can only help “closing the gap” for limited parametrizations

(E)NMPC parameters θ ought to be selected to minimize

$$J(\theta) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \ell(\mathbf{x}_k, \pi_\theta(\mathbf{x}_k)) \right] \quad \text{where } \mathbf{x}_{1,\dots,\infty} \text{ result from } \pi_\theta \rightarrow \text{real system}$$

rather than via model fitting!!

How to do that?

Reinforcement Learning - Core principles

Form function approximators:

$$Q_{\theta}(x, u), \quad V_{\theta}(x), \quad \pi_{\theta}(x)$$

via ad-hoc parametrization

Reinforcement Learning - Core principles

Form function approximators:

$$Q_\theta(x, u), V_\theta(x), \pi_\theta(x)$$

via ad-hoc parametrization

- **Q -learning methods** adjust θ to get

$$Q_*(x, u) \approx Q_\theta(x, u)$$

Yields policy:

$$\pi_\theta(x) = \underset{u}{\operatorname{arg\,min}} Q_\theta(x, u) \approx \underset{u}{\operatorname{arg\,min}} Q_*(x, u) = \pi_*(x)$$

Reinforcement Learning - Core principles

Form function approximators:

$$Q_\theta(x, u), V_\theta(x), \pi_\theta(x)$$

via ad-hoc parametrization

- **Q -learning methods** adjust θ to get

$$Q_*(x, u) \approx Q_\theta(x, u)$$

Yields policy:

$$\pi_\theta(x) = \underset{u}{\operatorname{arg\,min}} Q_\theta(x, u) \approx \underset{u}{\operatorname{arg\,min}} Q_*(x, u) = \pi_*(x)$$

- **Policy gradient methods** adjust θ to get

$$\nabla_\theta J(\pi_\theta) = 0$$

yields policy $\pi_\theta(x) \approx \pi_*(x)$ directly

Reinforcement Learning - Core principles

Form function approximators:

$$Q_\theta(x, u), V_\theta(x), \pi_\theta(x)$$

via ad-hoc parametrization

- **Q -learning methods** adjust θ to get

$$Q_*(x, u) \approx Q_\theta(x, u)$$

Yields policy:

$$\pi_\theta(x) = \underset{u}{\operatorname{arg\,min}} Q_\theta(x, u) \approx \underset{u}{\operatorname{arg\,min}} Q_*(x, u) = \pi_*(x)$$

- **Policy gradient methods** adjust θ to get

$$\nabla_\theta J(\pi_\theta) = 0$$

yields policy $\pi_\theta(x) \approx \pi_*(x)$ directly

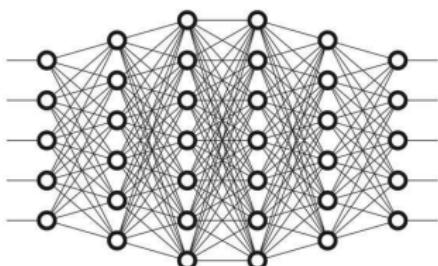
All approaches hinge on building either Q_θ or $\{\pi_\theta, V_\theta\}$

Reinforcement Learning - Core principles

Form function approximators:

$$Q_\theta(x, u), V_\theta(x), \pi_\theta(x)$$

via ad-hoc parametrization



- **Q -learning methods** adjust θ to get

$$Q_*(x, u) \approx Q_\theta(x, u)$$

Yields policy:

$$\pi_\theta(x) = \arg\min_u Q_\theta(x, u) \approx \arg\min_u Q_*(x, u) = \pi_*(x)$$

- **Policy gradient methods** adjust θ to get

$$\nabla_\theta J(\pi_\theta) = 0$$

yields policy $\pi_\theta(x) \approx \pi_*(x)$ directly

All approaches hinge on building either Q_θ or $\{\pi_\theta, V_\theta\}$

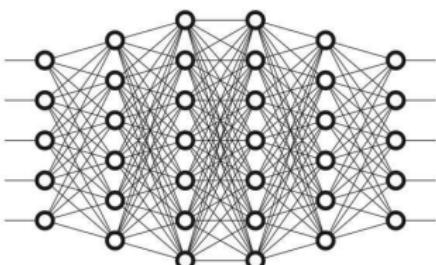
RL typically relies on DNNs as function approximators

Reinforcement Learning - Core principles

Form function approximators:

$$Q_\theta(x, u), V_\theta(x), \pi_\theta(x)$$

via ad-hoc parametrization



- **Q -learning methods** adjust θ to get

$$Q_*(x, u) \approx Q_\theta(x, u)$$

Yields policy:

$$\pi_\theta(x) = \arg \min_u Q_\theta(x, u) \approx \arg \min_u Q_*(x, u) = \pi_*(x)$$

- **Policy gradient methods** adjust θ to get

$$\nabla_\theta J(\pi_\theta) = 0$$

yields policy $\pi_\theta(x) \approx \pi_*(x)$ directly



All approaches hinge on building either Q_θ or $\{\pi_\theta, V_\theta\}$

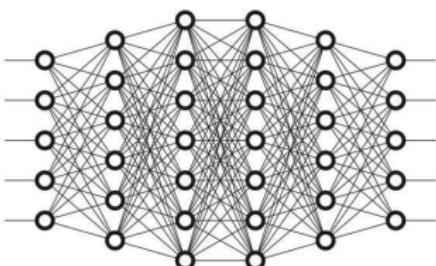
RL typically relies on DNNs as function approximators
i.e. beyond core principles, the rest is alchemy...

Reinforcement Learning - Core principles

Form function approximators:

$$Q_\theta(x, u), V_\theta(x), \pi_\theta(x)$$

via ad-hoc parametrization



- **Q -learning methods** adjust θ to get

$$Q_*(x, u) \approx Q_\theta(x, u)$$

Yields policy:

$$\pi_\theta(x) = \arg\min_u Q_\theta(x, u) \approx \arg\min_u Q_*(x, u) = \pi_*(x)$$

- **Policy gradient methods** adjust θ to get

$$\nabla_\theta J(\pi_\theta) = 0$$

yields policy $\pi_\theta(x) \approx \pi_*(x)$ directly



All approaches hinge on building either Q_θ or $\{\pi_\theta, V_\theta\}$

RL typically relies on DNNs as function approximators
i.e. beyond core principles, the rest is alchemy...

Can we be more structured and less obscure?

NMPC as a function approximator

NMPC delivers approximations of V and π

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

$$\pi_{\theta}(\mathbf{x}) = \mathbf{u}_0^* \quad (\text{ENMPC policy})$$

NMPC as a function approximator

NMPC delivers approximations of V and π

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

$$\pi_{\theta}(\mathbf{x}) = \mathbf{u}_0^* \quad (\text{ENMPC policy})$$

NMPC with constrained input delivers approx. of Q

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}, \quad \mathbf{u}_0 = \mathbf{u}$$

NMPC as a function approximator

NMPC delivers approximations of V and π

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}$$

$$\pi_{\theta}(\mathbf{x}) = \mathbf{u}_0^* \quad (\text{ENMPC policy})$$

NMPC with constrained input delivers approx. of Q

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{x}_0 = \mathbf{x}, \quad \mathbf{u}_0 = \mathbf{u}$$

DP relationships hold:

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}, \mathbf{u})$$

$$\pi_{\theta}(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}, \mathbf{u})$$

NMPC as a function approximator

NMPC delivers approximations of V and π

$$\begin{aligned} V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{x}} \quad & \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \quad & \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ & \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0 \\ & \mathbf{x}_0 = \mathbf{x} \end{aligned}$$

$$\pi_{\theta}(\mathbf{x}) = \mathbf{u}_0^* \quad (\text{ENMPC policy})$$

NMPC with constrained input delivers approx. of Q

$$\begin{aligned} Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{u}, \mathbf{x}} \quad & \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t. } \quad & \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ & \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0 \\ & \mathbf{x}_0 = \mathbf{x}, \quad \mathbf{u}_0 = \mathbf{u} \end{aligned}$$

DP relationships hold:

$$\begin{aligned} V_{\theta}(\mathbf{x}) &= \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}, \mathbf{u}) \\ \pi_{\theta}(\mathbf{x}) &= \arg \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}, \mathbf{u}) \end{aligned}$$

Theorem says this is a valid function approximator for RL. Bonus: strong theoretical properties

NMPC as a function approximator

NMPC delivers approximations of V and π

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$
 $\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$
 $\mathbf{x}_0 = \mathbf{x}$

$$\pi_{\theta}(\mathbf{x}) = \mathbf{u}_0^* \quad (\text{ENMPC policy})$$

NMPC with constrained input delivers approx. of Q

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$
 $\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$
 $\mathbf{x}_0 = \mathbf{x}, \quad \mathbf{u}_0 = \mathbf{u}$

DP relationships hold:

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}, \mathbf{u})$$
$$\pi_{\theta}(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_{\theta}(\mathbf{x}, \mathbf{u})$$

Theorem says this is a valid function approximator for RL. Bonus: strong theoretical properties

Use RL(-like) techniques to adjust the parameters θ

Outline

- 1 Introduction
- 2 Modification of the ENMPC scheme
- 3 Learning the Optimal ENMPC
- 4 RL for NMPC in practice
- 5 Simple example
- 6 Conclusions

Implementation

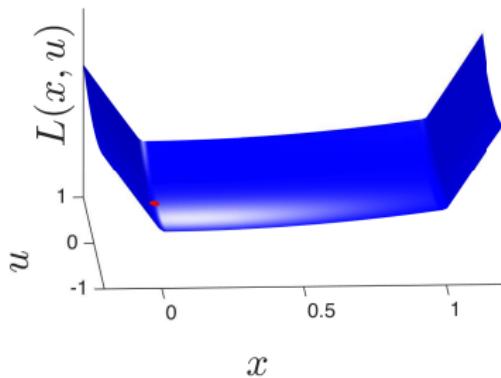
"Learning is impossible if mistakes yield an infinite punishment"

(E)NMPC scheme with mixed-constraints relaxation

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{x}, \mathbf{u}, \mathbf{s} \geq 0} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \left(L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{w}^\top \mathbf{s}_k \right)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$

$$\mathbf{h}_{\theta}^{\mathbf{u}}(\mathbf{u}_k) \leq 0, \quad \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq \mathbf{s}_k$$



'Walls' = relaxation of $x \in [0, 1]$

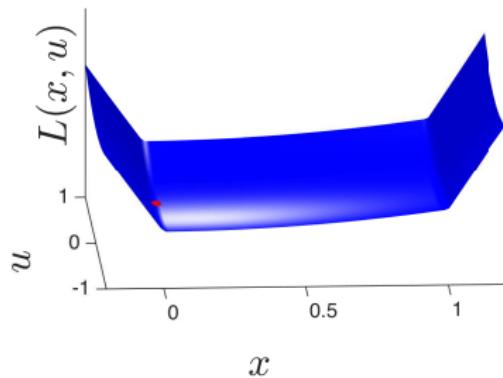
Implementation

"Learning is impossible if mistakes yield an infinite punishment"

(E)NMPC scheme with mixed-constraints relaxation

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{x}, \mathbf{u}, \mathbf{s} \geq 0} \overbrace{\lambda_{\theta}(\mathbf{x})}^{\rightarrow \text{dissipativity}} + \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \left(L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{w}^{\top} \mathbf{s}_k \right)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$
 $\mathbf{h}_{\theta}^{\mathbf{u}}(\mathbf{u}_k) \leq 0, \quad \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq \mathbf{s}_k$



'Walls' = relaxation of $x \in [0, 1]$

Implementation

"Learning is impossible if mistakes yield an infinite punishment"

(E)NMPC scheme with mixed-constraints relaxation

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{x}, \mathbf{u}, \mathbf{s} \geq 0} \overbrace{\lambda_{\theta}(\mathbf{x})}^{\rightarrow \text{dissipativity}} + \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \left(L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{w}^{\top} \mathbf{s}_k \right)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$
 $\mathbf{h}_{\theta}^{\mathbf{u}}(\mathbf{u}_k) \leq 0, \quad \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq \mathbf{s}_k$

TD error

$$\delta_{\theta}(\mathbf{x}, \mathbf{u}, \mathbf{x}_+) = L(\mathbf{x}, \mathbf{u}) + \gamma V_{\theta}(\mathbf{x}_+) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

E.g. use **LSTDQ**: solve

$$\mathbb{E}_{\tau_{\pi_{\theta}}} [\delta \nabla_{\theta} Q_{\theta}] = 0$$

for θ . This is on-the-fly LSPI for ENMPC

Implementation

"Learning is impossible if mistakes yield an infinite punishment"

(E)NMPC scheme with mixed-constraints relaxation

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{x}, \mathbf{u}, \mathbf{s} \geq 0} \overbrace{\lambda_{\theta}(\mathbf{x})}^{\rightarrow \text{dissipativity}} + \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \left(L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{w}^{\top} \mathbf{s}_k \right)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$
 $\mathbf{h}_{\theta}^{\mathbf{u}}(\mathbf{u}_k) \leq 0, \quad \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq \mathbf{s}_k$

TD error

$$\delta_{\theta}(\mathbf{x}, \mathbf{u}, \mathbf{x}_+) = L(\mathbf{x}, \mathbf{u}) + \gamma V_{\theta}(\mathbf{x}_+) - Q_{\theta}(\mathbf{x}, \mathbf{u})$$

E.g. use **LSTDQ**: solve

$$\mathbb{E}_{\tau_{\pi_{\theta}}} [\delta \nabla_{\theta} Q_{\theta}] = 0$$

for θ . This is on-the-fly LSPI for ENMPC

Can learn off-policy

- Generate new (E)NMPC θ' while using current one θ
- Update parameters $\theta' \leftarrow \theta$ after enough data and formal verifications

Implementation

"Learning is impossible if mistakes yield an infinite punishment"

(E)NMPC scheme with mixed-constraints relaxation

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{x}, \mathbf{u}, \mathbf{s} \geq 0} \overbrace{\lambda_{\theta}(\mathbf{x})}^{\rightarrow \text{dissipativity}} + \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \left(L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{w}^{\top} \mathbf{s}_k \right)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$
 $\mathbf{h}_{\theta}^{\mathbf{u}}(\mathbf{u}_k) \leq 0, \quad \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq \mathbf{s}_k$

Deterministic policy gradient

$$\min_{\theta} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \ell(\mathbf{x}_k, \pi_{\theta}(\mathbf{x}_k)) \right]$$

has necessary condition of optimality:

$$\nabla_{\theta} J = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta}(\mathbf{x}))] = 0$$

Implementation

"Learning is impossible if mistakes yield an infinite punishment"

(E)NMPC scheme with mixed-constraints relaxation

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{x}, \mathbf{u}, \mathbf{s} \geq 0} \overbrace{\lambda_{\theta}(\mathbf{x})}^{\rightarrow \text{dissipativity}} + \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \left(L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{w}^{\top} \mathbf{s}_k \right)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$
 $\mathbf{h}_{\theta}^{\mathbf{u}}(\mathbf{u}_k) \leq 0, \quad \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq \mathbf{s}_k$

Deterministic policy gradient

$$\min_{\theta} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \ell(\mathbf{x}_k, \pi_{\theta}(\mathbf{x}_k)) \right]$$

- Formal condition of optimality for tuned (E)NMPC scheme
- More complex to use than Q-learning

has necessary condition of optimality:

$$\nabla_{\theta} J = \mathbb{E}_{\tau_{\pi_{\theta}}} [\nabla_{\theta} \pi_{\theta}(\mathbf{x}) \nabla_{\mathbf{u}} Q_{\pi_{\theta}}(\mathbf{x}, \pi_{\theta}(\mathbf{x}))] = 0$$

Implementation

"Learning is impossible if mistakes yield an infinite punishment"

(E)NMPC scheme with mixed-constraints relaxation

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{x}, \mathbf{u}, \mathbf{s} \geq 0} \overbrace{\lambda_{\theta}(\mathbf{x})}^{\rightarrow \text{dissipativity}} + \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k \left(L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \mathbf{w}^{\top} \mathbf{s}_k \right)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$
 $\mathbf{h}_{\theta}^{\mathbf{u}}(\mathbf{u}_k) \leq 0, \quad \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq \mathbf{s}_k$

Theorem: if real dynamics + stage cost L are dissipative, then

- can use positive-definite L_{θ} and correction λ_{θ} such that ENMPC yields π_*, V_*, Q_* .
- λ_{θ} is dissipativity function
- ENMPC scheme is nominal-stable by construction

(valid for deterministic dynamics, should be extended)

Sensitivities

RL methods (may) require $\nabla_{\theta}\pi_{\theta}$, $\nabla_u Q_{\theta}$, $\nabla_{\theta}Q_{\theta}$, $\nabla_{\theta}^2Q_{\theta}$, etc. how to get that?

Sensitivities

RL methods (may) require $\nabla_{\theta}\pi_{\theta}$, $\nabla_{\mathbf{u}}Q_{\theta}$, $\nabla_{\theta}Q_{\theta}$, $\nabla_{\theta}^2Q_{\theta}$, etc. how to get that?

NMPC delivers approx. of Q

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{u}_0 = \mathbf{u}$$

Sensitivities

RL methods (may) require $\nabla_{\theta}\pi_{\theta}$, $\nabla_{\mathbf{u}}Q_{\theta}$, $\nabla_{\theta}Q_{\theta}$, $\nabla_{\theta}^2Q_{\theta}$, etc. how to get that?

NMPC delivers approx. of Q

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{u}_0 = \mathbf{u}$$

Consider the Lagrange function:

$$\begin{aligned} \mathcal{L} = & \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \sum_{k=1}^N \boldsymbol{\lambda}_k^\top (\mathbf{x}_{k+1} - \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)) + \sum_{k=0}^{N-1} \boldsymbol{\mu}_k^\top \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ & + \boldsymbol{\lambda}_0^\top (\mathbf{x}_0 - \mathbf{x}) + \boldsymbol{\nu}^\top (\mathbf{u}_0 - \mathbf{u}) \end{aligned}$$

Sensitivities

RL methods (may) require $\nabla_{\theta} \pi_{\theta}$, $\nabla_{\mathbf{u}} Q_{\theta}$, $\nabla_{\theta} Q_{\theta}$, $\nabla_{\theta}^2 Q_{\theta}$, etc. how to get that?

NMPC delivers approx. of Q

$$Q_{\theta}(\mathbf{x}, \mathbf{u}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

$$\mathbf{u}_0 = \mathbf{u}$$

Consider the Lagrange function:

$$\begin{aligned} \mathcal{L} = & \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \sum_{k=1}^N \boldsymbol{\lambda}_k^\top (\mathbf{x}_{k+1} - \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)) + \sum_{k=0}^{N-1} \boldsymbol{\mu}_k^\top \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ & + \boldsymbol{\lambda}_0^\top (\mathbf{x}_0 - \mathbf{x}) + \boldsymbol{\nu}^\top (\mathbf{u}_0 - \mathbf{u}) \end{aligned}$$

Then

$$\nabla_{\theta} Q_{\theta}(\mathbf{x}, \mathbf{u}) = \nabla_{\theta} \mathcal{L}, \quad \nabla_{\mathbf{u}} Q_{\theta} = \nabla_{\mathbf{u}} \mathcal{L}$$

... evaluated on the primal-dual solution of the ENMPC scheme

Sensitivities

RL methods (may) require $\nabla_{\theta}\pi_{\theta}$, $\nabla_{\mathbf{u}}Q_{\theta}$, $\nabla_{\theta}Q_{\theta}$, $\nabla_{\theta}^2Q_{\theta}$, etc. how to get that?

NMPC delivers approx. of V

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

Sensitivities

RL methods (may) require $\nabla_{\theta}\pi_{\theta}$, $\nabla_{\mathbf{u}}Q_{\theta}$, $\nabla_{\theta}Q_{\theta}$, $\nabla_{\theta}^2Q_{\theta}$, etc. how to get that?

NMPC delivers approx. of V

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

Consider the Lagrange function:

$$\begin{aligned} \mathcal{L} = & \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \sum_{k=1}^N \boldsymbol{\lambda}_k^\top (\mathbf{x}_{k+1} - \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)) + \sum_{k=0}^{N-1} \boldsymbol{\mu}_k^\top \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ & + \boldsymbol{\lambda}_0^\top (\mathbf{x}_0 - \mathbf{x}) \end{aligned}$$

Sensitivities

RL methods (may) require $\nabla_{\theta}\pi_{\theta}$, $\nabla_{\mathbf{u}}Q_{\theta}$, $\nabla_{\theta}Q_{\theta}$, $\nabla_{\theta}^2Q_{\theta}$, etc. how to get that?

NMPC delivers approx. of V

$$V_{\theta}(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

Consider the Lagrange function:

$$\mathcal{L} = \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) + \sum_{k=1}^N \boldsymbol{\lambda}_k^\top (\mathbf{x}_{k+1} - \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)) + \sum_{k=0}^{N-1} \boldsymbol{\mu}_k^\top \mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$
$$+ \boldsymbol{\lambda}_0^\top (\mathbf{x}_0 - \mathbf{x})$$

Then

$$\nabla_{\theta} V_{\theta}(\mathbf{x}) = \nabla_{\theta} \mathcal{L}$$

... evaluated on the primal-dual solution of the ENMPC scheme

Sensitivities

RL methods (may) require $\nabla_{\theta}\pi_{\theta}$, $\nabla_{\mathbf{u}}Q_{\theta}$, $\nabla_{\theta}Q_{\theta}$, $\nabla_{\theta}^2Q_{\theta}$, etc. how to get that?

NMPC delivers policy π_{θ}

$$\pi_{\theta}(\mathbf{x}) = \mathbf{u}_0^*$$

$$\mathbf{u}^*, \mathbf{x}^* = \arg \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k), \quad \mathbf{x}_0 = \mathbf{x}$$

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

Sensitivity $\nabla_{\theta}\pi_{\theta}(\mathbf{x}) = \nabla_{\theta}\mathbf{u}_0^*$, computed from:

$$\frac{\partial \mathbf{z}}{\partial \theta} = -\frac{\partial \mathbf{R}}{\partial \mathbf{z}}^{-1} \frac{\partial \mathbf{R}}{\partial \theta}$$

Second-order sensitivities e.g.

$$\nabla_{\theta}^2 V_{\theta} = \frac{\partial^2 \mathcal{L}}{\partial \theta^2} + \frac{\partial^2 \mathcal{L}}{\partial \theta \partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \theta}$$

where

- \mathbf{R} are the KKT conditions of ENMPC
- \mathbf{z} collects the primal-dual variables
- Best treated in an primal-dual interior-point framework with "large" barrier parameter

Outline

- 1 Introduction
- 2 Modification of the ENMPC scheme
- 3 Learning the Optimal ENMPC
- 4 RL for NMPC in practice
- 5 Simple example
- 6 Conclusions

Q-learning on our trivial example

Modified NMPC scheme

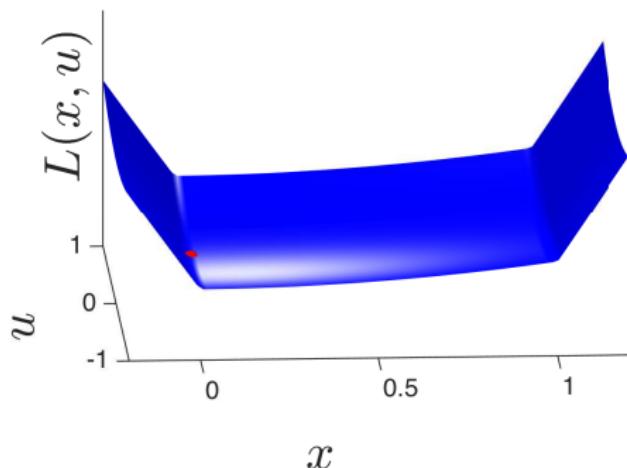
$$\begin{aligned}\pi(x) \leftarrow \min_{u, x, s} \quad & \gamma^N V_*(x_N) + \sum_{k=0}^{N-1} \gamma^k \hat{L}(x, u) \\ \text{s.t.} \quad & x_{k+1} = x_k + 0.1u_k \\ & -1 \leq u_k \leq 1 \\ & x_0 = x\end{aligned}$$

Real dynamics:

$$x_+ = x + 0.1u + e$$

Noise:

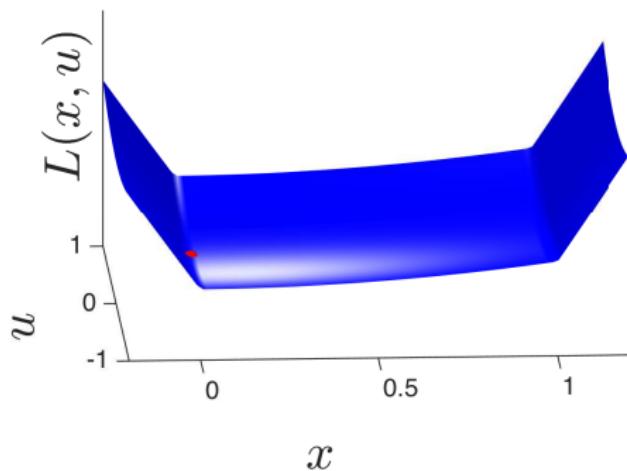
$$e \sim \mathcal{U}([-0.1, 0])$$



Q -learning on our trivial example

Modified NMPC scheme

$$\begin{aligned}\pi(x) \leftarrow \min_{u, x, s} \quad & \gamma^N V_*(x_N) + \sum_{k=0}^{N-1} \gamma^k \hat{L}(x, u) \\ \text{s.t.} \quad & x_{k+1} = x_k + 0.1u_k \\ & -1 \leq u_k \leq 1 \\ & x_0 = x\end{aligned}$$

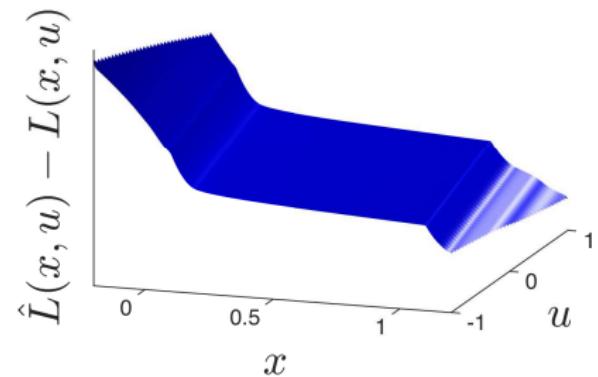


Real dynamics:

$$x_+ = x + 0.1u + e$$

Noise:

$$e \sim \mathcal{U}([-0.1, 0])$$



Q-learning on our trivial example

NMPC with relaxed constraints (\equiv "walls")

$$\min_{u, x, s} \quad c + \gamma^N T x_N^2 + \sum_{k=0}^{N-1} \gamma^k \left(x_k^2 + u_k^2 + ws_k + ws_k^2 + f^\top \begin{bmatrix} x_k \\ u_k \end{bmatrix} \right)$$

$$\text{s.t. } x_{k+1} = x_k + 0.1u_k + b$$

$$-1 \leq u_k \leq 1, \quad \underline{x} - s_k \leq x_k \leq \bar{x} + s_k$$

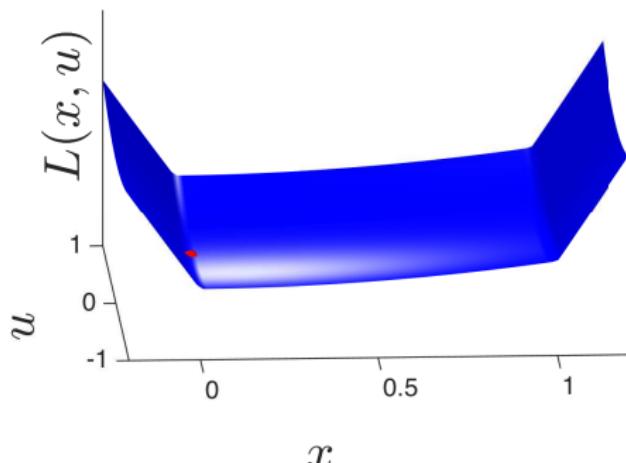
$$x_0 = x$$

Real dynamics:

$$x_+ = x + 0.1u + e$$

Noise:

$$e \sim \mathcal{U}([-0.1, 0])$$



Q-learning can tune all parameters:

- Model bias b
- Cost gradient f
- State constraints (walls position) \underline{x}, \bar{x}
- Terminal cost T

Q -learning on our trivial example

NMPC with relaxed constraints (\equiv "walls")

$$\min_{u, x, s} \quad \textcolor{blue}{c} + \gamma^N \textcolor{blue}{T} x_N^2 + \sum_{k=0}^{N-1} \gamma^k \left(x_k^2 + u_k^2 + ws_k + ws_k^2 + \mathbf{f}^\top \begin{bmatrix} x_k \\ u_k \end{bmatrix} \right)$$

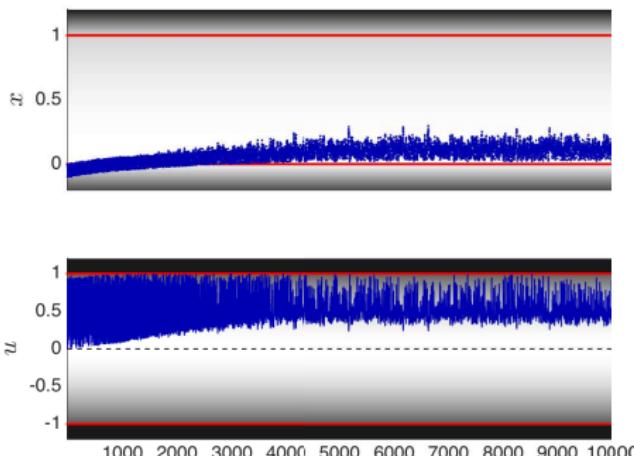
s.t. $x_{k+1} = x_k + 0.1u_k + \textcolor{blue}{b}$
 $-1 \leq u_k \leq 1, \quad \underline{x} - s_k \leq x_k \leq \bar{x} + s_k$
 $x_0 = x$

Real dynamics:

$$x_+ = x + 0.1u + e$$

Noise:

$$e \sim \mathcal{U}([-0.1, 0])$$



Q -learning on our trivial example

NMPC with relaxed constraints (\equiv "walls")

$$\min_{u, x, s} \quad \mathbf{c} + \gamma^N \mathbf{T} x_N^2 + \sum_{k=0}^{N-1} \gamma^k \left(x_k^2 + u_k^2 + ws_k + ws_k^2 + \mathbf{f}^\top \begin{bmatrix} x_k \\ u_k \end{bmatrix} \right)$$

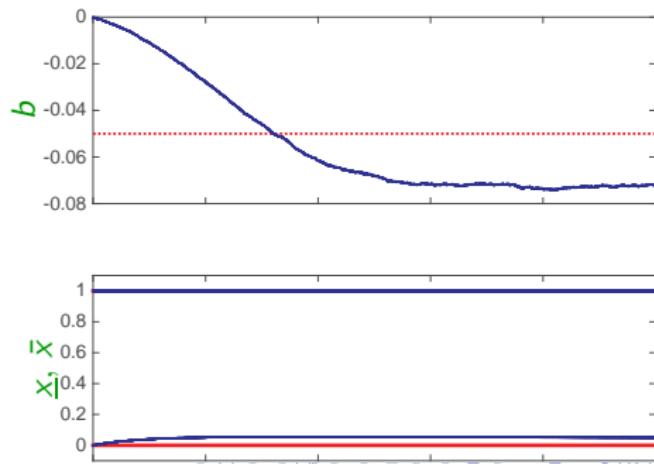
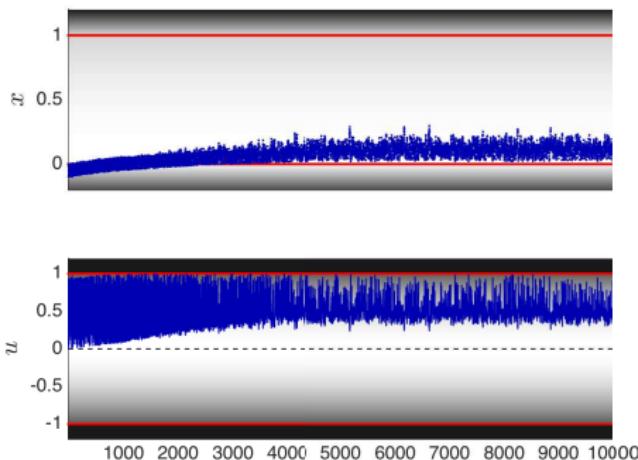
s.t. $x_{k+1} = x_k + 0.1u_k + b$
 $-1 \leq u_k \leq 1, \quad \underline{x} - s_k \leq x_k \leq \bar{x} + s_k$
 $x_0 = x$

Real dynamics:

$$x_+ = x + 0.1u + e$$

Noise:

$$e \sim \mathcal{U}([-0.1, 0])$$



Q-learning on our trivial example

NMPC with relaxed constraints (\equiv "walls")

$$\min_{u, x, s} \quad \textcolor{blue}{c} + \gamma^N \textcolor{blue}{T} x_N^2 + \sum_{k=0}^{N-1} \gamma^k \left(x_k^2 + u_k^2 + ws_k + ws_k^2 + \mathbf{f}^\top \begin{bmatrix} x_k \\ u_k \end{bmatrix} \right)$$

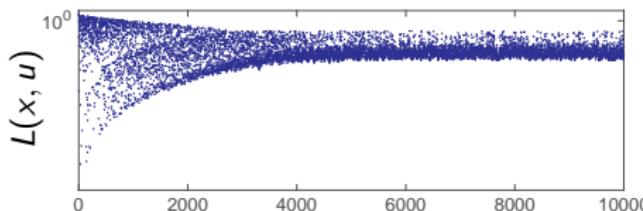
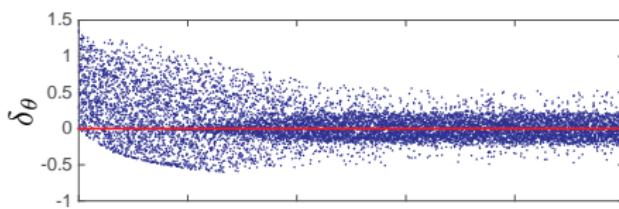
s.t. $x_{k+1} = x_k + 0.1u_k + \textcolor{blue}{b}$
 $-1 \leq u_k \leq 1, \quad \textcolor{blue}{x} - s_k \leq x_k \leq \bar{x} + s_k$
 $x_0 = x$

Real dynamics:

$$x_+ = x + 0.1u + e$$

Noise:

$$e \sim \mathcal{U}([-0.1, 0])$$



Q -learning on our trivial example

NMPC with relaxed constraints (\equiv "walls")

$$\min_{u, x, s} \quad \textcolor{blue}{c} + \gamma^N \textcolor{blue}{T} x_N^2 + \sum_{k=0}^{N-1} \gamma^k \left(x_k^2 + u_k^2 + ws_k + ws_k^2 + \mathbf{f}^\top \begin{bmatrix} x_k \\ u_k \end{bmatrix} \right)$$

s.t. $x_{k+1} = x_k + 0.1u_k + \textcolor{blue}{b}$

$$-1 \leq u_k \leq 1, \quad \underline{x} - s_k \leq x_k \leq \bar{x} + s_k$$

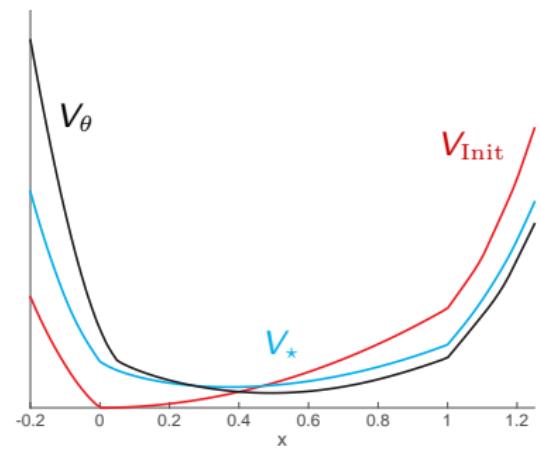
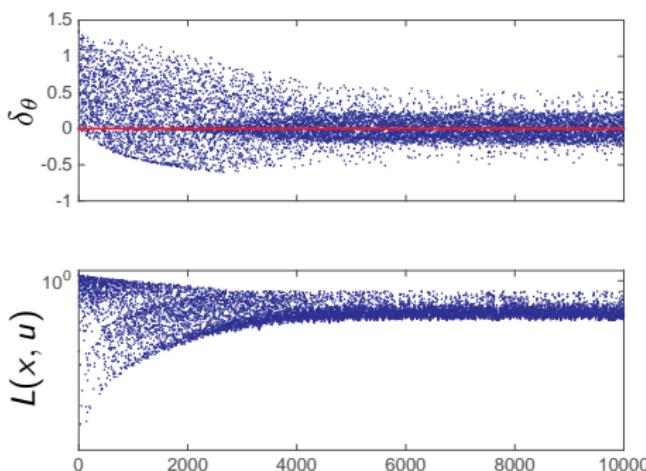
$$x_0 = x$$

Real dynamics:

$$x_+ = x + 0.1u + e$$

Noise:

$$e \sim \mathcal{U}([-0.1, 0])$$



Q -learning on our trivial example

NMPC with relaxed constraints (\equiv “walls”)

$$\min_{u, x, s} \quad \textcolor{blue}{c} + \gamma^N \textcolor{blue}{T} x_N^2 + \sum_{k=0}^{N-1} \gamma^k \left(x_k^2 + u_k^2 + ws_k + ws_k^2 + \mathbf{f}^\top \begin{bmatrix} x_k \\ u_k \end{bmatrix} \right)$$

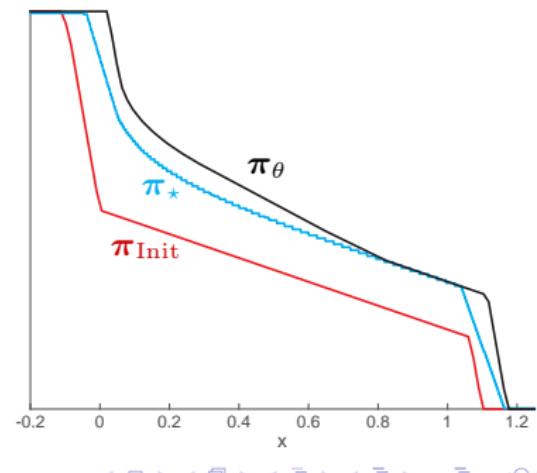
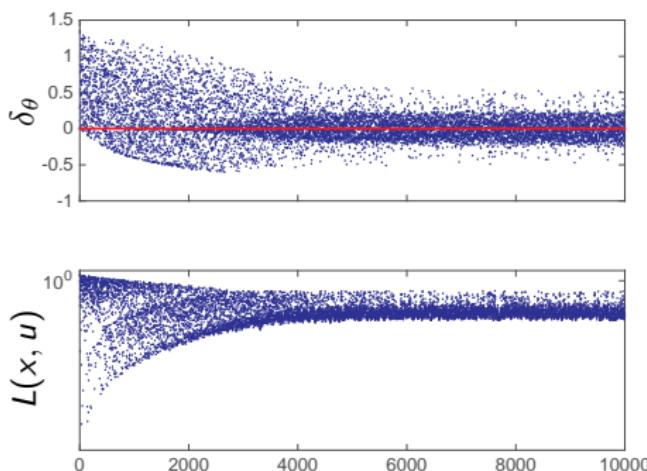
$$\begin{aligned} \text{s.t. } \quad & x_{k+1} = x_k + 0.1u_k + \textcolor{blue}{b} \\ & -1 \leq u_k \leq 1, \quad \textcolor{blue}{x} - s_k \leq x_k \leq \bar{x} + s_k \\ & x_0 = x \end{aligned}$$

Real dynamics:

$$x_+ = x + 0.1u + e$$

Noise:

$$e \sim \mathcal{U}([-0.1, 0])$$



Outline

- 1 Introduction
- 2 Modification of the ENMPC scheme
- 3 Learning the Optimal ENMPC
- 4 RL for NMPC in practice
- 5 Simple example
- 6 Conclusions

Where is SYSID then?

Tune parameters such that $\nabla_{\theta} J(\pi_{\theta}) = 0$. Model fitting does not matter?

Where is SYSID then?

Tune parameters such that $\nabla_{\theta} J(\pi_{\theta}) = 0$. Model fitting does not matter?

- Model fitting the data is probably a good place to start with

Where is SYSID then?

Tune parameters such that $\nabla_{\theta} J(\pi_{\theta}) = 0$. Model fitting does not matter?

- Model fitting the data is probably a good place to start with
- $\nabla_{\theta}^2 Q_{\theta}$ and $\nabla_{\theta}^2 J(\pi_{\theta})$ are (typically) rank-deficient. Some parameters are left to be adjusted → SYSID as “regularization”

Where is SYSID then?

Tune parameters such that $\nabla_{\theta} J(\pi_{\theta}) = 0$. Model fitting does not matter?

- Model fitting the data is probably a good place to start with
- $\nabla_{\theta}^2 Q_{\theta}$ and $\nabla_{\theta}^2 J(\pi_{\theta})$ are (typically) rank-deficient. Some parameters are left to be adjusted → SYSID as “regularization”
- Safe learning...

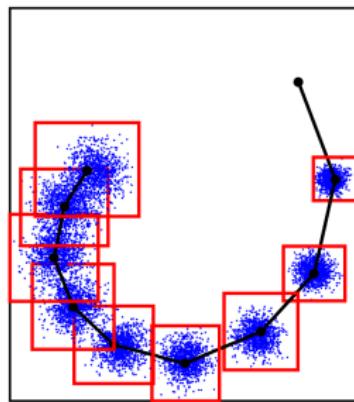
Where is SYSID then?

Tune parameters such that $\nabla_{\theta} J(\pi_{\theta}) = 0$. Model fitting does not matter?

- Model fitting the data is probably a good place to start with
- $\nabla_{\theta}^2 Q_{\theta}$ and $\nabla_{\theta}^2 J(\pi_{\theta})$ are (typically) rank-deficient. Some parameters are left to be adjusted → SYSID as “regularization”
- Safe learning...

Deterministic “forward-set” model:

$$\mathbf{x}, \mathbf{u} \rightarrow \mathbf{x}_+ \in X_{\theta}^{+}(\mathbf{x}, \mathbf{u}, \theta)$$

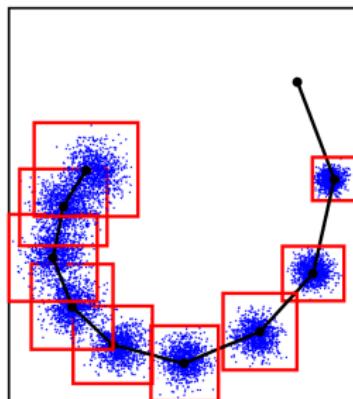


Where is SYSID then?

Tune parameters such that $\nabla_{\theta} J(\pi_{\theta}) = 0$. Model fitting does not matter?

- Model fitting the data is probably a good place to start with
- $\nabla_{\theta}^2 Q_{\theta}$ and $\nabla_{\theta}^2 J(\pi_{\theta})$ are (typically) rank-deficient. Some parameters are left to be adjusted → SYSID as “regularization”
- Safe learning...

Deterministic “forward-set” model:



$$\mathbf{x}, \mathbf{u} \rightarrow \mathbf{x}_+ \in X_{\theta}^{+}(\mathbf{x}, \mathbf{u}, \theta)$$

Robust ENMPC optimizes under constraint:

$$\begin{aligned} \min_{\mathbf{u}, \mathbf{x}} \quad & \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ \text{s.t.} \quad & X_{\theta}^0 = \mathbf{x}, \quad X_{\theta}^k \subseteq \mathbb{X} \end{aligned}$$

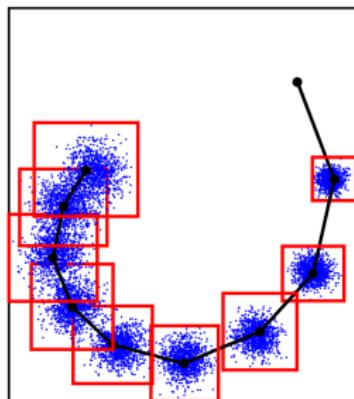
where X_{θ}^k is the set-propagation of the uncertainties.

Where is SYSID then?

Tune parameters such that $\nabla_{\theta} J(\pi_{\theta}) = 0$. Model fitting does not matter?

- Model fitting the data is probably a good place to start with
- $\nabla_{\theta}^2 Q_{\theta}$ and $\nabla_{\theta}^2 J(\pi_{\theta})$ are (typically) rank-deficient. Some parameters are left to be adjusted → SYSID as “regularization”
- Safe learning...

Deterministic “forward-set” model:



$$\mathbf{x}, \mathbf{u} \rightarrow \mathbf{x}_+ \in X_{\theta}^{+}(\mathbf{x}, \mathbf{u}, \theta)$$

Robust ENMPC optimizes under constraint:

$$\begin{aligned} & \min_{\mathbf{u}, \mathbf{x}} \gamma^N T_{\theta}(\mathbf{x}_N) + \sum_{k=0}^{N-1} \gamma^k L_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \\ & \text{s.t. } X_{\theta}^0 = \mathbf{x}, \quad X_{\theta}^k \subseteq \mathbb{X} \end{aligned}$$

where X_{θ}^k is the set-propagation of the uncertainties.

Adjust θ to ensure X_{θ}^{+} “encloses the data”, and minimize $J(\pi_{\theta})$

Conclusions & Future work

Conclusions

- Cost function (and constraints) of an ENMPC scheme ought not necessarily match the ones describing the problem, should be adjusted

Data-driven Economic NMPC using Reinforcement Learning, S. Gros, M. Zanon, Transaction on Automatic Control
A Data-Efficient Framework for Learning NMPC controllers, M. Zanon, S. Gros, A. Bemporad, ECC 2019 Deterministic Policy Gradient for Data-Driven Economic NMPC, S. Gros, M. Zanon
System Identification and Reinforcement Learning for Data-Driven Economic NMPC, S. Gros, M. Zanon

Conclusions & Future work

Conclusions

- Cost function (and constraints) of an ENMPC scheme ought not necessarily match the ones describing the problem, should be adjusted
- Adjusting ENMPC scheme for performance \neq model fitting

Data-driven Economic NMPC using Reinforcement Learning, S. Gros, M. Zanon, Transaction on Automatic Control
A Data-Efficient Framework for Learning NMPC controllers, M. Zanon, S. Gros, A. Bemporad, ECC 2019 Deterministic Policy Gradient for Data-Driven Economic NMPC, S. Gros, M. Zanon
System Identification and Reinforcement Learning for Data-Driven Economic NMPC, S. Gros, M. Zanon

Conclusions & Future work

Conclusions

- Cost function (and constraints) of an ENMPC scheme ought not necessarily match the ones describing the problem, should be adjusted
- Adjusting ENMPC scheme for performance \neq model fitting
- Necessary Conditions of Optimality of an ENMPC scheme are defined by policy gradient theorems. Can be deployed to tune ENMPC schemes.

Data-driven Economic NMPC using Reinforcement Learning, S. Gros, M. Zanon, Transaction on Automatic Control
A Data-Efficient Framework for Learning NMPC controllers, M. Zanon, S. Gros, A. Bemporad, ECC 2019 Deterministic Policy Gradient for Data-Driven Economic NMPC, S. Gros, M. Zanon
System Identification and Reinforcement Learning for Data-Driven Economic NMPC, S. Gros, M. Zanon

Conclusions & Future work

Conclusions

- Cost function (and constraints) of an ENMPC scheme ought not necessarily match the ones describing the problem, should be adjusted
- Adjusting ENMPC scheme for performance \neq model fitting
- Necessary Conditions of Optimality of an ENMPC scheme are defined by policy gradient theorems. Can be deployed to tune ENMPC schemes.
- ENMPC schemes can be used as a structured function approximation for Reinforcement Learning, more insights than DNN, opens new possibilities

Data-driven Economic NMPC using Reinforcement Learning, S. Gros, M. Zanon, Transaction on Automatic Control
A Data-Efficient Framework for Learning NMPC controllers, M. Zanon, S. Gros, A. Bemporad, ECC 2019 Deterministic Policy Gradient for Data-Driven Economic NMPC, S. Gros, M. Zanon
System Identification and Reinforcement Learning for Data-Driven Economic NMPC, S. Gros, M. Zanon

Conclusions & Future work

Conclusions

- Cost function (and constraints) of an ENMPC scheme ought not necessarily match the ones describing the problem, should be adjusted
- Adjusting ENMPC scheme for performance \neq model fitting
- Necessary Conditions of Optimality of an ENMPC scheme are defined by policy gradient theorems. Can be deployed to tune ENMPC schemes.
- ENMPC schemes can be used as a structured function approximation for Reinforcement Learning, more insights than DNN, opens new possibilities
- Strong connections to dissipativity theory in ENMPC (not covered today)

Data-driven Economic NMPC using Reinforcement Learning, S. Gros, M. Zanon, Transaction on Automatic Control
A Data-Efficient Framework for Learning NMPC controllers, M. Zanon, S. Gros, A. Bemporad, ECC 2019 Deterministic Policy Gradient for Data-Driven Economic NMPC, S. Gros, M. Zanon
System Identification and Reinforcement Learning for Data-Driven Economic NMPC, S. Gros, M. Zanon

Conclusions & Future work

Conclusions

- Cost function (and constraints) of an ENMPC scheme ought not necessarily match the ones describing the problem, should be adjusted
- Adjusting ENMPC scheme for performance \neq model fitting
- Necessary Conditions of Optimality of an ENMPC scheme are defined by policy gradient theorems. Can be deployed to tune ENMPC schemes.
- ENMPC schemes can be used as a structured function approximation for Reinforcement Learning, more insights than DNN, opens new possibilities
- Strong connections to dissipativity theory in ENMPC (not covered today)

Future work

- Computationally efficient RL methods for ENMPC (LSTD/experience replay)

Data-driven Economic NMPC using Reinforcement Learning, S. Gros, M. Zanon, Transaction on Automatic Control
A Data-Efficient Framework for Learning NMPC controllers, M. Zanon, S. Gros, A. Bemporad, ECC 2019
Deterministic Policy Gradient for Data-Driven Economic NMPC, S. Gros, M. Zanon
System Identification and Reinforcement Learning for Data-Driven Economic NMPC, S. Gros, M. Zanon

Conclusions & Future work

Conclusions

- Cost function (and constraints) of an ENMPC scheme ought not necessarily match the ones describing the problem, should be adjusted
- Adjusting ENMPC scheme for performance \neq model fitting
- Necessary Conditions of Optimality of an ENMPC scheme are defined by policy gradient theorems. Can be deployed to tune ENMPC schemes.
- ENMPC schemes can be used as a structured function approximation for Reinforcement Learning, more insights than DNN, opens new possibilities
- Strong connections to dissipativity theory in ENMPC (not covered today)

Future work

- Computationally efficient RL methods for ENMPC (LSTD/experience replay)
- SYSID & RL for ENMPC context? It's a safe-learning question...

Data-driven Economic NMPC using Reinforcement Learning, S. Gros, M. Zanon, Transaction on Automatic Control
A Data-Efficient Framework for Learning NMPC controllers, M. Zanon, S. Gros, A. Bemporad, ECC 2019
Deterministic Policy Gradient for Data-Driven Economic NMPC, S. Gros, M. Zanon
System Identification and Reinforcement Learning for Data-Driven Economic NMPC, S. Gros, M. Zanon

Conclusions & Future work

Conclusions

- Cost function (and constraints) of an ENMPC scheme ought not necessarily match the ones describing the problem, should be adjusted
- Adjusting ENMPC scheme for performance \neq model fitting
- Necessary Conditions of Optimality of an ENMPC scheme are defined by policy gradient theorems. Can be deployed to tune ENMPC schemes.
- ENMPC schemes can be used as a structured function approximation for Reinforcement Learning, more insights than DNN, opens new possibilities
- Strong connections to dissipativity theory in ENMPC (not covered today)

Future work

- Computationally efficient RL methods for ENMPC (LSTD/experience replay)
- SYSID & RL for ENMPC context? It's a safe-learning question...
- Big Data for ENMPC tuning?

Data-driven Economic NMPC using Reinforcement Learning, S. Gros, M. Zanon, Transaction on Automatic Control
A Data-Efficient Framework for Learning NMPC controllers, M. Zanon, S. Gros, A. Bemporad, ECC 2019
Deterministic Policy Gradient for Data-Driven Economic NMPC, S. Gros, M. Zanon
System Identification and Reinforcement Learning for Data-Driven Economic NMPC, S. Gros, M. Zanon