

Learning + MPC via Reinforcement Learning

Fundamental principles

Sébastien Gros

Dept. of Cybernetic, NTNU
Faculty of Information Tech.

Freiburg PhD School

On this topic

- First publication in 2020
- ~40 papers
- Many talks & courses
- Growing portfolio of applications & experiments
- A bit on the “theoretical” side in the field

On these lectures

- Give high-level concepts
- Focus on known insights
- What are the current gaps
- New insights (3rd lecture)

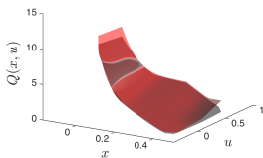
Software for implementation are not mature yet. You will be the first “large” audience playing with them.



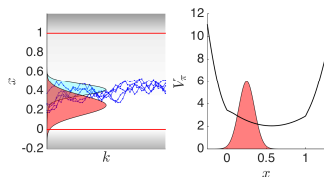
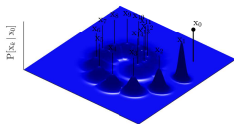
What are we going to discuss?

- 1 Learning for MPC - A focus on closed-loop performance
- 2 Safety & stability in Learning for MPC
- 3 When do “classic” approaches work / When is learning beneficial?

samples = 1000000



$$Q_+(x, u) \leftarrow L(x, u) + \gamma \mathbb{E} [V(x_+) \mid x, u]$$



Outline

- 1 The Basics
- 2 More background
- 3 Let's take a deeper dive
- 4 Parametrization & Role of the model
- 5 RL over MPC

Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state s

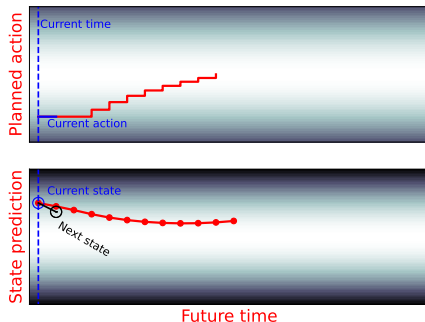
$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state s

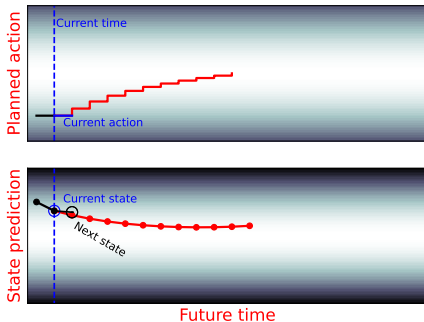
$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state \mathbf{s}

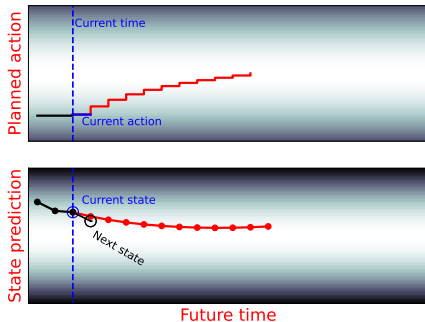
$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state \mathbf{s}

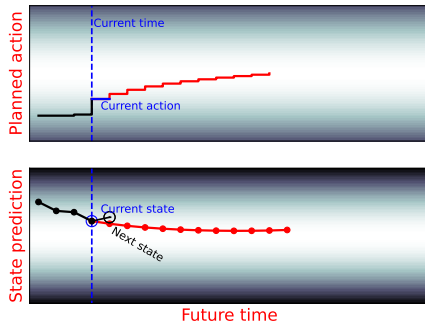
$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state \mathbf{s}

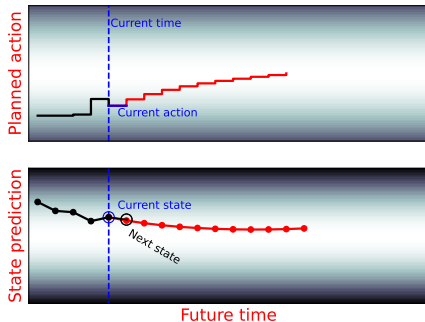
$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state \mathbf{s}

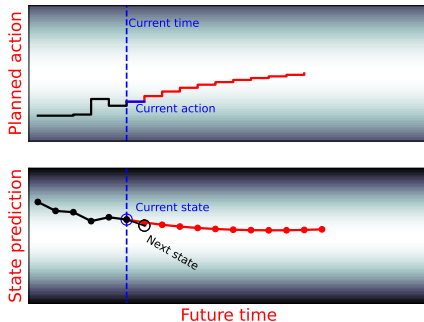
$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state s

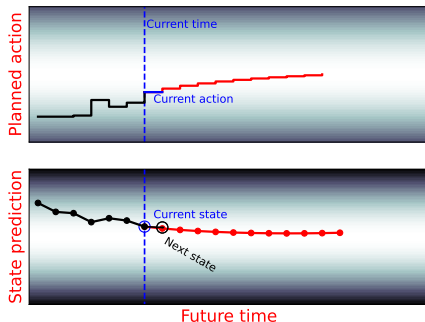
$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state s

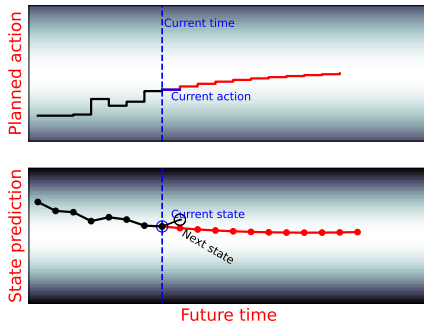
$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state s

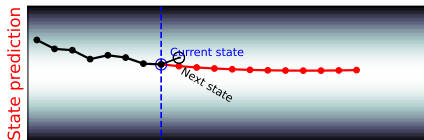
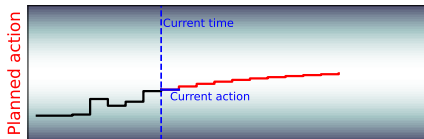
$$\min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



MPC

- is based on planning the future
- Policy from repeated planning

$$\pi^{\text{MPC}}(s) = \mathbf{u}_0^*$$

Model Predictive Control (MPC)

Optimize a plan over finite horizon, apply first move, repeat

MPC: at current state s

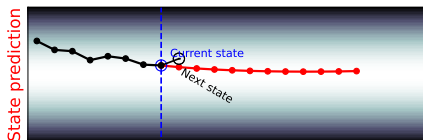
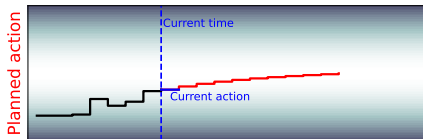
$$\min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

apply action $\mathbf{a} = \mathbf{u}_0^*$ to the system



MPC

- is based on planning the future
- Policy from repeated planning

$$\pi^{\text{MPC}}(s) = \mathbf{u}_0^*$$

MPC is a powerful tool to control constrained systems, increasingly used as a practical way of building optimal policies

Theoretical Framework to connect RL and MPC

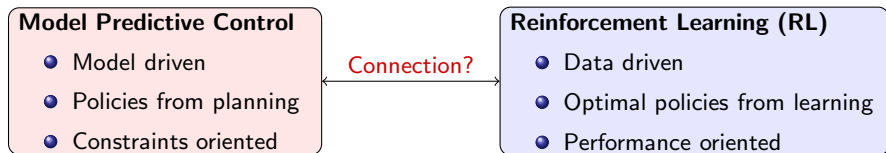
Model Predictive Control

- Model driven
- Policies from planning
- Constraints oriented

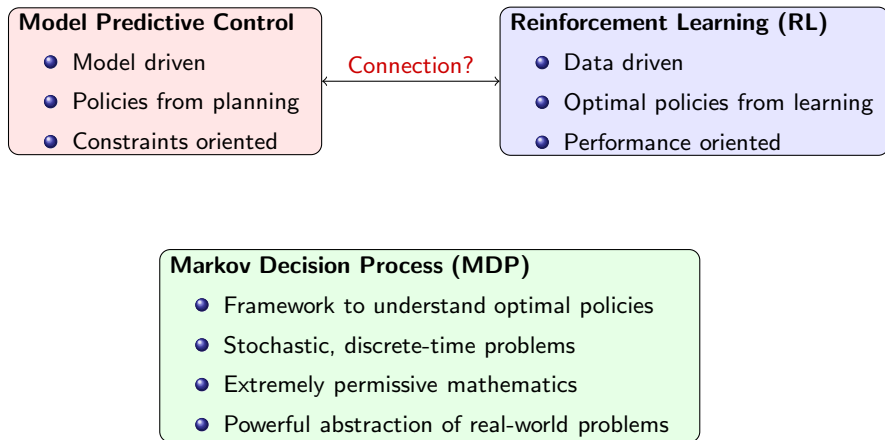
Reinforcement Learning (RL)

- Data driven
- Optimal policies from learning
- Performance oriented

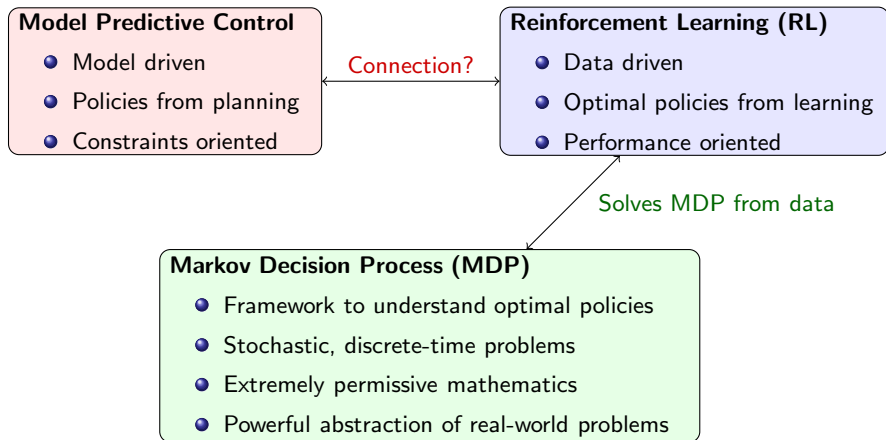
Theoretical Framework to connect RL and MPC



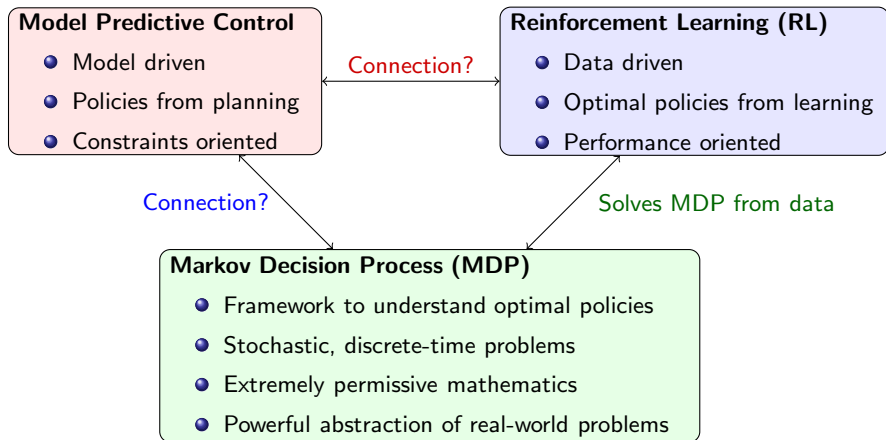
Theoretical Framework to connect RL and MPC



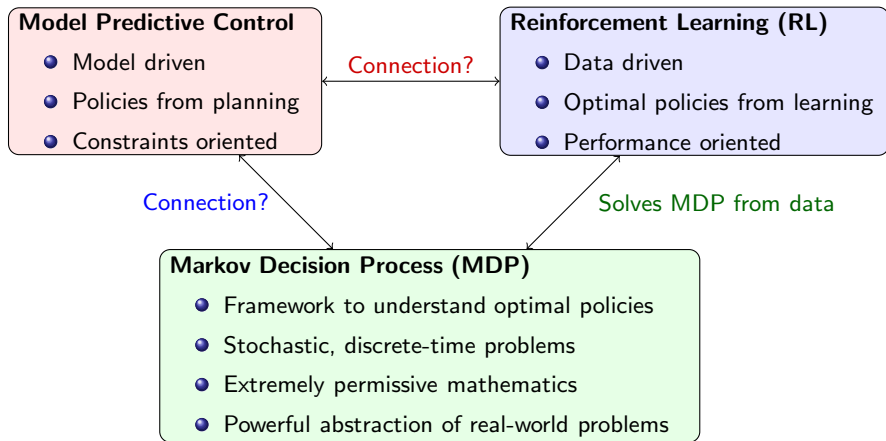
Theoretical Framework to connect RL and MPC



Theoretical Framework to connect RL and MPC



Theoretical Framework to connect RL and MPC



Connecting MPC and RL is about connecting MPC to MDPs!!

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, \mathbf{a} \rightarrow s_+$$

(state-action \rightarrow next state)

Cost function (instant performance)

$$L(s, \mathbf{a}) \in \mathbb{R}$$

A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$\mathbf{s}, \mathbf{a} \rightarrow \mathbf{s}_+$$

(state-action \rightarrow next state)

Cost function (instant performance)

$$L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$$

- **Policy**

$$\mathbf{a} = \pi(\mathbf{s})$$

is how we act on the system

A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$\mathbf{s}, \mathbf{a} \rightarrow \mathbf{s}_+$$

(state-action \rightarrow next state)

Cost function (instant performance)

$$L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$$

- **Policy**

$$\mathbf{a} = \boldsymbol{\pi}(\mathbf{s})$$

is how we act on the system

- **Closed-loop performance**

$$J(\boldsymbol{\pi}) = \mathbb{E}_{\boldsymbol{\pi}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \boldsymbol{\pi}(\mathbf{s}_k)) \right]$$

with discount $\gamma \in [0, 1]$

A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

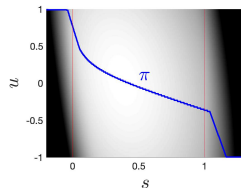
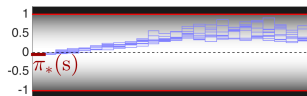
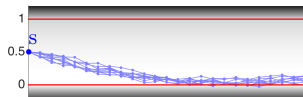
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

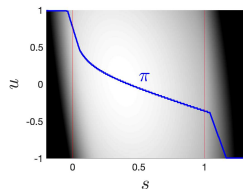
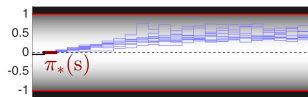
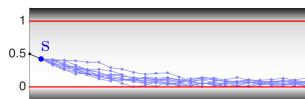
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

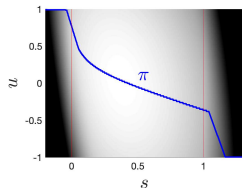
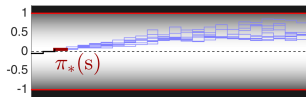
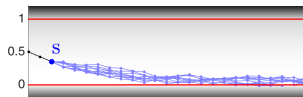
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

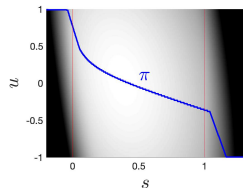
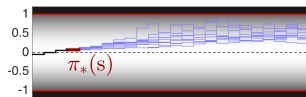
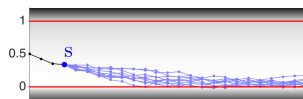
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

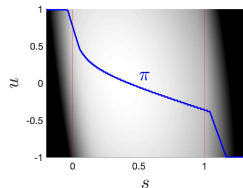
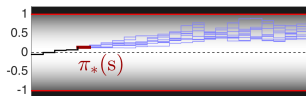
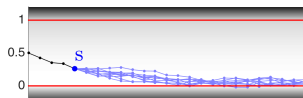
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- Policy

$$a = \pi(s)$$

is how we act on the system

- Closed-loop performance

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

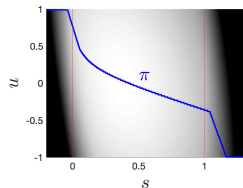
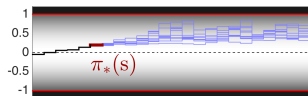
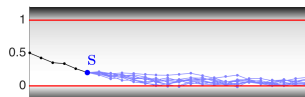
with discount $\gamma \in [0, 1]$

- Optimal policy: π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

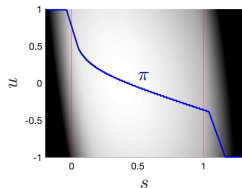
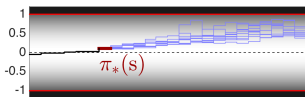
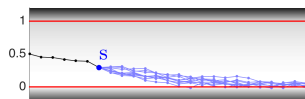
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

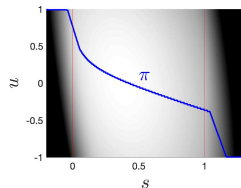
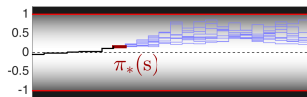
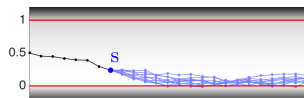
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

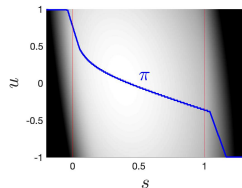
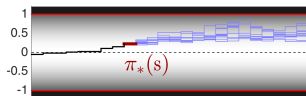
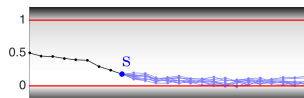
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

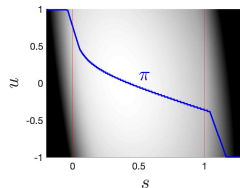
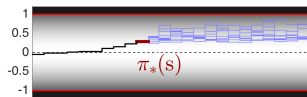
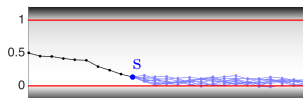
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

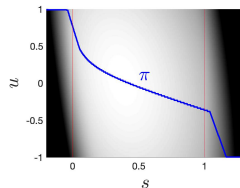
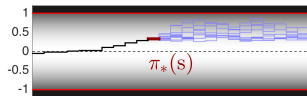
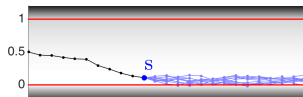
with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$



A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$

MDP is a go-to framework when considering general optimal control problems, useful for applications with stochastic dynamics.

A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$s, a \rightarrow s_+$$

(state-action \rightarrow next state)

- **Policy**

$$a = \pi(s)$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, \pi(s_k)) \right]$$

with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(s, a) \in \mathbb{R}$$

MDP is a go-to framework when considering general optimal control problems, useful for applications with stochastic dynamics.

Solution of an MDP is described by “simple” equations, but solving them is very challenging

A (fairly) general way of describing optimal control

Markov Decision Processes (MDP)

Stochastic state transitions

$$\mathbf{s}, \mathbf{a} \rightarrow \mathbf{s}_+$$

(state-action \rightarrow next state)

- **Policy**

$$\mathbf{a} = \pi(\mathbf{s})$$

is how we act on the system

- **Closed-loop performance**

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \pi(\mathbf{s}_k)) \right]$$

with discount $\gamma \in [0, 1]$

- **Optimal policy:** π^* from

$$\min_{\pi} J(\pi)$$

Cost function (instant performance)

$$L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$$

MDP is a go-to framework when considering general optimal control problems, useful for applications with stochastic dynamics.

Solution of an MDP is described by “simple” equations, but solving them is very challenging

By doing “re-planning” all the time, MPC generates a policy π^{MPC} that *hopefully* resembles π^*

MPC is a heuristic to solve MDPs

why do we use it?

A (fairly) general way of describing optimal control

Purpose of MPC? *According to the MPC community*

Historically MPC focuses on **constraints satisfaction & stability**, track a reference
Tracking MPC

More recent focus is on **closed-loop performance**, e.g. energy, time, money.
Economic MPC

Purpose of MPC? *According to the MPC community*

Historically MPC focuses on **constraints satisfaction & stability**, track a reference
Tracking MPC

More recent focus is on **closed-loop performance**, e.g. energy, time, money.
Economic MPC

E.g. of the form:

$$\min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}^T W \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}$$

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- Costs are “designed” to steer the system to reference “(0, 0)”
- MPC is not optimizing a specific “physical quantity” (cost unit??)

Purpose of MPC? According to the MPC community

Historically MPC focuses on **constraints satisfaction & stability**, track a reference
Tracking MPC

More recent focus is on **closed-loop performance**, e.g. energy, time, money.
Economic MPC

E.g. of the form:

$$\min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}^T W \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}$$

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- Costs are “designed” to steer the system to reference “(0, 0)”
- MPC is not optimizing a specific “physical quantity” (cost unit??)

Generic form:

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

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- L directly represents something we want to minimize
- MPC is optimizing a specific “physical quantity” (cost has a unit)

Purpose of MPC? According to the MPC community

Historically MPC focuses on **constraints satisfaction & stability**, track a reference
Tracking MPC

More recent focus is on **closed-loop performance**, e.g. energy, time, money.
Economic MPC

E.g. of the form:

$$\min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}^T W \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}$$

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- Costs are “designed” to steer the system to reference “(0, 0)”
- MPC is not optimizing a specific “physical quantity” (cost unit??)

Generic form:

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

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- L directly represents something we want to minimize
- MPC is optimizing a specific “physical quantity” (cost has a unit)

“MPC is for constraints satisfaction...” (heard in scientific discussions)

Purpose of MPC? According to the MPC community

Historically MPC focuses on **constraints satisfaction & stability**, track a reference
Tracking MPC

More recent focus is on **closed-loop performance**, e.g. energy, time, money.
Economic MPC

E.g. of the form:

$$\min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}^T W \begin{bmatrix} \mathbf{x}_k \\ \mathbf{u}_k \end{bmatrix}$$

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- Costs are “designed” to steer the system to reference “(0, 0)”
- MPC is not optimizing a specific “physical quantity” (cost unit??)

Generic form:

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

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- L directly represents something we want to minimize
- MPC is optimizing a specific “physical quantity” (cost has a unit)

“MPC is for constraints satisfaction...” (heard in scientific discussions)

... it is, but it does not need to be limited to that.

Optimal policies from MPC?

Optimality often cast as minimizing[†]

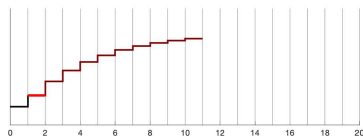
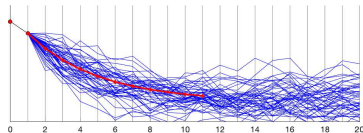
$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$



Optimal policies from MPC?

Optimality often cast as minimizing[†]

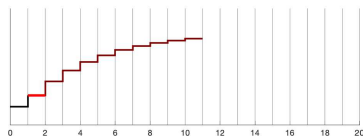
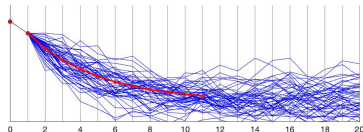
$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$



Optimal policies from MPC?

Optimality often cast as minimizing[†]

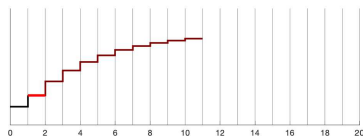
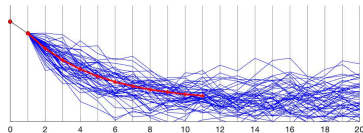
$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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

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

$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

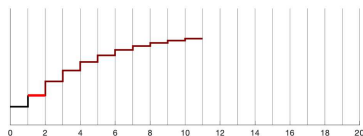
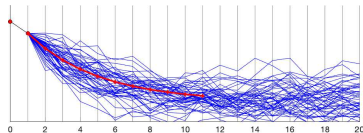
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

[†]Alternative forms of optimality: bias / gain optimal, beyond 1st moment

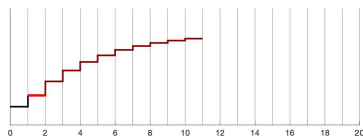
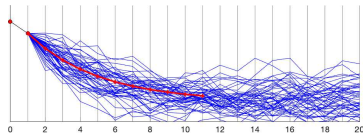
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

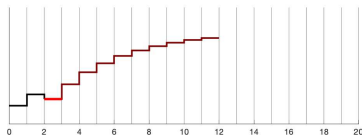
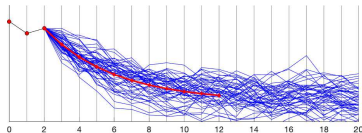
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

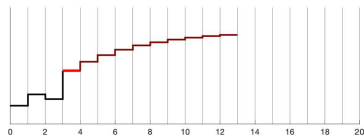
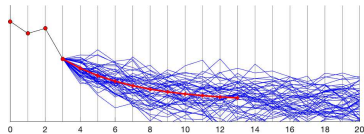
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

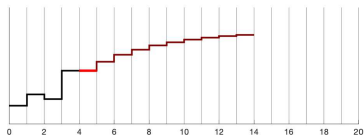
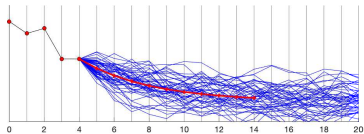
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

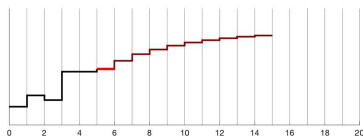
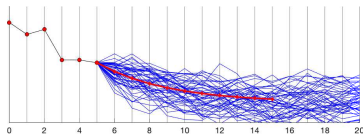
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

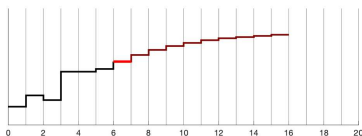
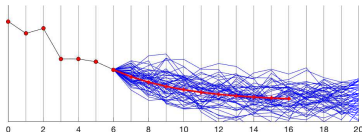
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

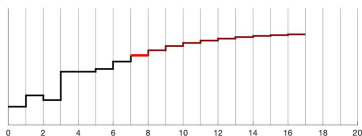
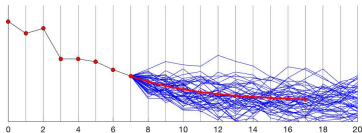
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model **f inaccurate**

... in general no

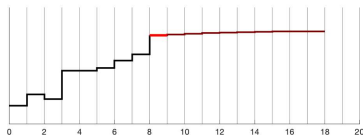
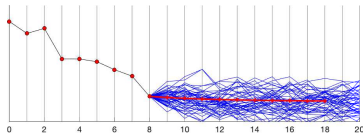
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

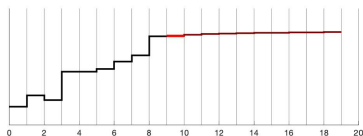
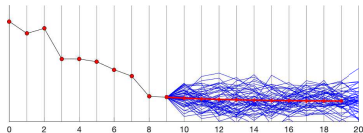
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

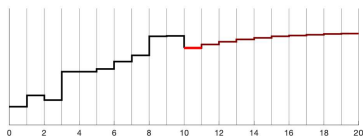
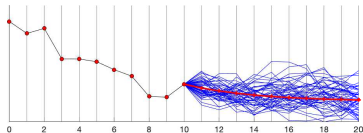
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

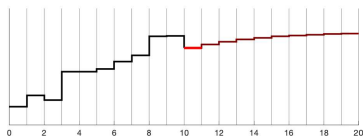
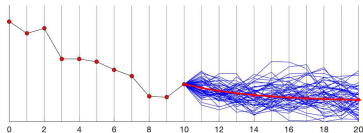
Optimal policies from MPC?

Optimality often cast as minimizing[†]

$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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



Optimal policy π^* :

$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

Optimal policies from MPC?

Optimality often cast as minimizing[†]

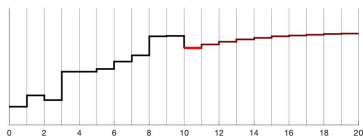
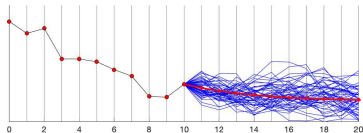
$$J(\pi) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

MPC policy $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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

Can learning help with that?

...different "bets"



Optimal policy π^* :

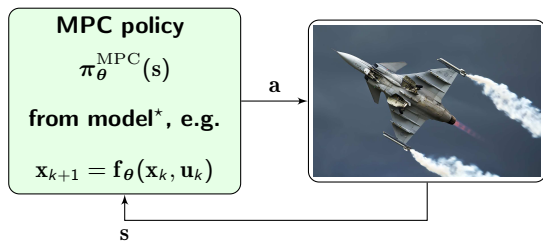
$$\pi^* = \underset{\pi}{\text{a min}} J(\pi)$$

Does MPC give $\pi^{\text{MPC}} = \pi^*$?

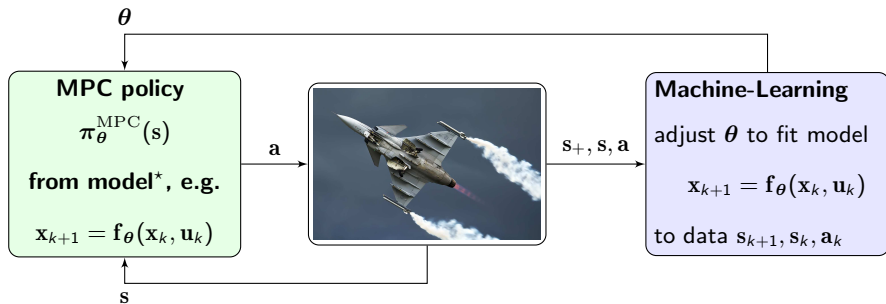
- 1 Infinite vs. finite horizon
- 2 Discounted ($\gamma < 1$) vs. undiscounted
- 3 Model \mathbf{f} **inaccurate**

... in general no

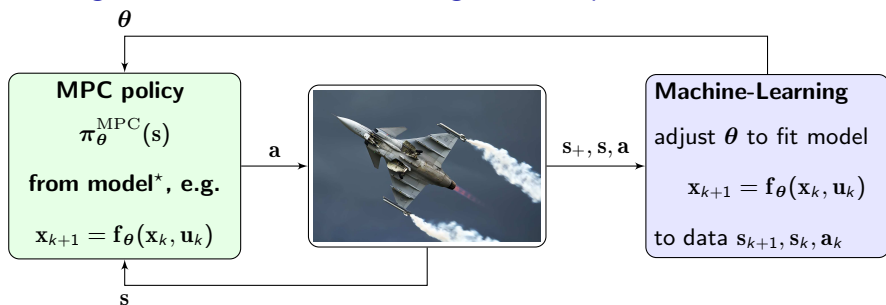
Learning for MPC - Machine Learning in-the-loop



Learning for MPC - Machine Learning in-the-loop



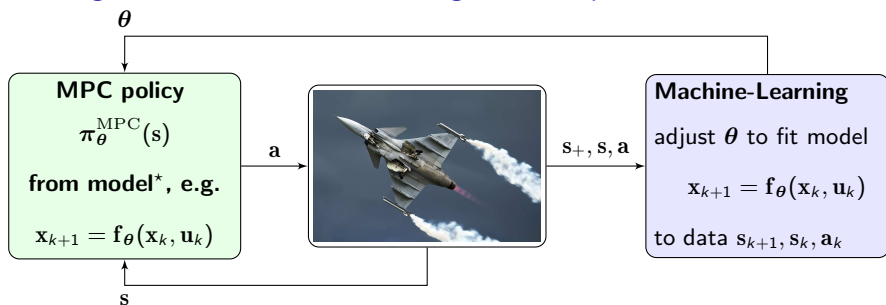
Learning for MPC - Machine Learning in-the-loop



“Machine-Learning” in-the-loop f_{θ} from

- **Physics-based:** first principles + SYSID
- **Neural Network:** DNN, LSTM, TFT, ...
- **Statistical:** GP, RKHS, GPC, ARX ...

Learning for MPC - Machine Learning in-the-loop

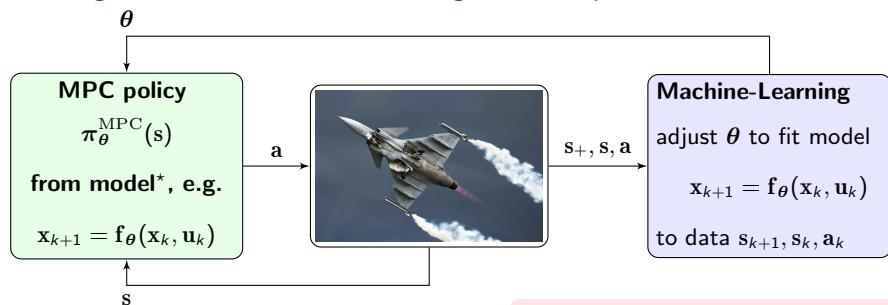


“Machine-Learning” in-the-loop \mathbf{f}_{θ} from

- **Physics-based:** first principles + SYSID
- **Neural Network:** DNN, LSTM, TFT, ...
- **Statistical:** GP, RKHS, GPC, ARX ...

** can replace “model” by any prediction strategies:
input-output predictors, multi-step predictors, etc...*

Learning for MPC - Machine Learning in-the-loop



“Machine-Learning” in-the-loop \mathbf{f}_{θ} from

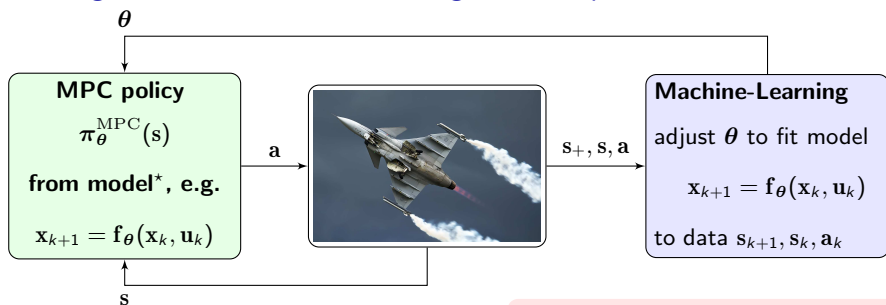
- **Physics-based:** first principles + SYSID
- **Neural Network:** DNN, LSTM, TFT, ...
- **Statistical:** GP, RKHS, GPC, ARX ...

Paradigm

- Performance tied to prediction accuracy
- Target accuracy via ML
- Ignore that MPC is a policy

*can replace “model” by any prediction strategies:
input-output predictors, multi-step predictors, etc...

Learning for MPC - Machine Learning in-the-loop



“Machine-Learning” in-the-loop \mathbf{f}_{θ} from

- **Physics-based:** first principles + SYSID
- **Neural Network:** DNN, LSTM, TFT, ...
- **Statistical:** GP, RKHS, GPC, ARX ...

* can replace “model” by any prediction strategies:
input-output predictors, multi-step predictors, etc...

Paradigm

- Performance tied to prediction accuracy
- Target accuracy via ML
- Ignore that MPC is a policy

We focus on “breaking” this paradigm
Learning / RL plays a key role

Paradigm shifts...

Paradigm shifts...

Shift 1: focus on performance instead of fitting

- **from:** f_θ is a model for the system dynamics
- **to:** MPC is a model of optimality (will specify that in a bit...)

Paradigm shifts...

Shift 1: focus on performance instead of fitting

- **from:** f_θ is a model for the system dynamics
- **to:** MPC is a model of optimality (will specify that in a bit...)

Classic view...

MPC: at current state \mathbf{s} solve

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

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

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

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

gives policy $\pi_\theta^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$

Find θ such that prediction
“fits” the data

Paradigm shifts...

Shift 1: focus on performance instead of fitting

- **from:** f_θ is a model for the system dynamics
- **to:** MPC is a model of optimality (will specify that in a bit...)

Classic view...

MPC: at current state \mathbf{s} solve

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

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

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

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

gives policy $\pi_\theta^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$

Shift to...

Find θ that “fits MPC to optimality”
according to the data, e.g. minimizes $J(\pi_\theta^{\text{MPC}})$

Find θ such that prediction
“fits” the data

Paradigm shifts...

Shift 1: focus on performance instead of fitting

- **from:** f_θ is a model for the system dynamics
- **to:** MPC is a model of optimality (will specify that in a bit...)

Classic view...

MPC: at current state \mathbf{s} solve

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

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

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

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

gives policy $\pi_\theta^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$

Shift to...

Find θ that “fits MPC to optimality”
according to the data, e.g. minimizes $J(\pi_\theta^{\text{MPC}})$

- \rightarrow Best model for closed-loop performance
- \neq Best model to fit the data!
- *More on this in 3rd lecture*

RL is a toolbox to do that...

Find θ such that prediction
“fits” the data

Paradigm shifts...

Shift 1: focus on performance instead of fitting

- **from:** \mathbf{f}_θ is a model for the system dynamics
- **to:** MPC is a model of optimality (will specify that in a bit...)

Classic view...

MPC: at current state \mathbf{s} solve

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

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

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

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

gives policy $\pi_\theta^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$

Find θ such that prediction
“fits” the data

Shift to...

Find θ that “fits MPC to optimality”
according to the data, e.g. minimizes $J(\pi_\theta^{\text{MPC}})$

- \rightarrow Best model for closed-loop performance
- \neq Best model to fit the data!
- *More on this in 3rd lecture*

RL is a toolbox to do that...

But getting π^ places “high demands” on \mathbf{f}_θ*

Paradigm shifts...

Shift 1: focus on performance instead of fitting

- **from:** f_θ is a model for the system dynamics
- **to:** MPC is a model of optimality (will specify that in a bit...)

Classic view...

MPC: at current state \mathbf{s} solve

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

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

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

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

gives policy $\pi_\theta^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$

Find θ such that prediction
“fits” the data

Shift to...

Find θ that “fits MPC to optimality”
according to the data, e.g. minimizes $J(\pi_\theta^{\text{MPC}})$

- \rightarrow Best model for closed-loop performance
- \neq Best model to fit the data!
- *More on this in 3rd lecture*

RL is a toolbox to do that...

But getting π^ places “high demands” on f_θ*

Can we do more? **Yes...**

Paradigm shifts...

Shift 1: focus on performance instead of fitting

- **from:** f_θ is a model for the system dynamics
- **to:** MPC is a model of optimality (will specify that in a bit...)

Classic view...

MPC: at current state \mathbf{s} solve

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

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

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

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

gives policy $\pi_\theta^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$

Find θ such that prediction
“fits” the data

Shift to...

Find θ that “fits MPC to optimality”
according to the data, e.g. minimizes $J(\pi_\theta^{\text{MPC}})$

Shift 2: “holistic” parametrization

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

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

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

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

gives policy $\pi_\theta^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$

How to use this? Reinforcement Learning

Policy $\pi_{\theta}^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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

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

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- $\min_{\theta} J(\pi_{\theta}^{\text{MPC}})$ using data
- $\theta \rightarrow J(\pi_{\theta}^{\text{MPC}})$ very implicit
- $J(\cdot)$ is the real-system!

How to use this? Reinforcement Learning

Reinforcement Learning

Tools to approximate π^* from data
This is not (necessarily) about DNNs

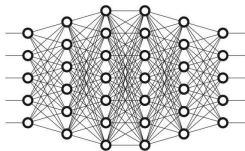
Policy $\pi_{\theta}^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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

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

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- $\min_{\theta} J(\pi_{\theta}^{\text{MPC}})$ using data
- $\theta \rightarrow J(\pi_{\theta}^{\text{MPC}})$ very implicit
- $J(\cdot)$ is the real-system!



How to use this? Reinforcement Learning

Reinforcement Learning

Tools to approximate π^* from data
This is not (necessarily) about DNNs

For MPC: tools to find best θ , e.g.

- **Policy Gradient:** estimations of $\nabla_{\theta} J(\pi_{\theta}^{\text{MPC}})$, possibly $\nabla_{\theta}^2 J(\pi_{\theta}^{\text{MPC}})$
- **Q-learning:** direct “shaping” of MPC

Combination is useful...

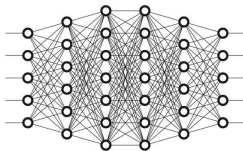
Policy $\pi_{\theta}^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ from

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

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

$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

- $\min_{\theta} J(\pi_{\theta}^{\text{MPC}})$ using data
- $\theta \rightarrow J(\pi_{\theta}^{\text{MPC}})$ very implicit
- $J(\cdot)$ is the real-system!



Outline

- 1 The Basics
- 2 More background**
- 3 Let's take a deeper dive
- 4 Parametrization & Role of the model
- 5 RL over MPC

Optimal Value Functions

- **Value function:**

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\boldsymbol{\pi}_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_k = \boldsymbol{\pi}_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy $\boldsymbol{\pi}_{\star}$, starting from given initial conditions \mathbf{s}

Optimal Value Functions

- **Value function:**

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_k = \pi_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy π_{\star} , starting from given initial conditions \mathbf{s}

- **Action-Value function:**

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}, \mathbf{a}_{k>0} = \pi_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy π_{\star} , starting from given initial conditions \mathbf{s} , and using action \mathbf{a} as first input (policy π_{\star} after that)

Optimal Value Functions

- **Value function:**

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_k = \pi_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy π_{\star} , starting from given initial conditions \mathbf{s}

- **Action-Value function:**

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}, \mathbf{a}_{k>0} = \pi_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy π_{\star} , starting from given initial conditions \mathbf{s} , and using action \mathbf{a} as first input (policy π_{\star} after that)

- **Relationship:**

$$V_{\star}(\mathbf{s}) = \min_{\mathbf{a}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

Optimal Value Functions

- **Value function:**

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_k = \pi_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy π_{\star} , starting from given initial conditions \mathbf{s}

- **Action-Value function:**

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}, \mathbf{a}_{k>0} = \pi_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy π_{\star} , starting from given initial conditions \mathbf{s} , and using action \mathbf{a} as first input (policy π_{\star} after that)

- **Relationship:**

$$V_{\star}(\mathbf{s}) = \min_{\mathbf{a}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

- **Optimal Policy:**

$$\pi_{\star}(\mathbf{s}) = \arg \min_{\mathbf{a}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

Optimal Value Functions

- **Value function:**

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_k = \pi_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy π_{\star} , starting from given initial conditions \mathbf{s}

- **Action-Value function:**

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}, \mathbf{a}_{k>0} = \pi_{\star}(\mathbf{s}_k) \right]$$

gives the expected cost for policy π_{\star} , starting from given initial conditions \mathbf{s} , and using action \mathbf{a} as first input (policy π_{\star} after that)

- **Relationship:**

$$V_{\star}(\mathbf{s}) = \min_{\mathbf{a}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

- **Optimal Policy:**

$$\pi_{\star}(\mathbf{s}) = \arg \min_{\mathbf{a}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

Can be computed via the Bellman equations, intractable for “large” state-action spaces

Value Functions

- **Value function:**

$$V_{\pi}(\mathbf{s}) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid s_0 = \mathbf{s}, \mathbf{a}_k = \pi(\mathbf{s}_k) \right]$$

gives the expected cost for policy π , starting from given initial conditions \mathbf{s}

- **Action-Value function:**

$$Q_{\pi}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid s_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a}, \mathbf{a}_{k>0} = \pi(\mathbf{s}_k) \right]$$

gives the expected cost for policy π , starting from given initial conditions \mathbf{s} , and using action \mathbf{a} as first input (policy π_{\star} after that)

- **Relationship:**

$$V_{\pi}(\mathbf{s}) = Q_{\pi}(\mathbf{s}, \pi(\mathbf{s}_k))$$

Note:

$$V_{\pi} \neq V_{\star}$$

- **Advantage function:**

$$Q_{\pi} \neq Q_{\star}$$

$$A_{\pi}(\mathbf{s}, \mathbf{a}) = Q_{\pi}(\mathbf{s}, \mathbf{a}) - V_{\pi}(\mathbf{s})$$

$$A_{\pi} \neq A_{\star}$$

compares \mathbf{a} to policy π . Instrumental in policy gradient methods.

Can be computed via the Bellman equations, intractable for “large” state-action

MDPs and “forbidden” states

What if the system is not allowed to leave a certain subset of the state space?

What if the system is not allowed to leave a certain subset of the state space?

- Say there is a “feasible” set:

$$\mathbb{F} = \{ \mathbf{s} \mid \mathbf{h}(\mathbf{s}) \leq 0 \}$$

where the state of the system should always be.

What if the system is not allowed to leave a certain subset of the state space?

- Say there is a “feasible” set:

$$\mathbb{F} = \{ \mathbf{s} \mid \mathbf{h}(\mathbf{s}) \leq 0 \}$$

where the state of the system should always be.

- In the “MDP theory”, assign an infinite penalty to leaving \mathbb{F} , i.e. add:

$$I_{\mathbb{F}}(\mathbf{s}, \mathbf{a}) = \begin{cases} 0 & \text{if } \mathbf{s} \in \mathbb{F} \\ +\infty & \text{if } \mathbf{s} \notin \mathbb{F} \end{cases}$$

to stage cost L .

What if the system is not allowed to leave a certain subset of the state space?

- Say there is a “feasible” set:

$$\mathbb{F} = \{ \mathbf{s} \mid \mathbf{h}(\mathbf{s}) \leq 0 \}$$

where the state of the system should always be.

- In the “MDP theory”, assign an infinite penalty to leaving \mathbb{F} , i.e. add:

$$I_{\mathbb{F}}(\mathbf{s}, \mathbf{a}) = \begin{cases} 0 & \text{if } \mathbf{s} \in \mathbb{F} \\ +\infty & \text{if } \mathbf{s} \notin \mathbb{F} \end{cases}$$

to stage cost L .

- In RL, ∞ penalties are not meaningful: “There is no backup from death”

What if the system is not allowed to leave a certain subset of the state space?

- Say there is a “feasible” set:

$$\mathbb{F} = \{ \mathbf{s} \mid \mathbf{h}(\mathbf{s}) \leq 0 \}$$

where the state of the system should always be.

- In the “MDP theory”, assign an infinite penalty to leaving \mathbb{F} , i.e. add:

$$I_{\mathbb{F}}(\mathbf{s}, \mathbf{a}) = \begin{cases} 0 & \text{if } \mathbf{s} \in \mathbb{F} \\ +\infty & \text{if } \mathbf{s} \notin \mathbb{F} \end{cases}$$

to stage cost L .

- In RL, ∞ penalties are not meaningful: “There is no backup from death”
- Common approach: assign a “very large” penalty to $\mathbf{s} \notin \mathbb{F}$ instead of $+\infty$.

What if the system is not allowed to leave a certain subset of the state space?

- Say there is a “feasible” set:

$$\mathbb{F} = \{ \mathbf{s} \mid \mathbf{h}(\mathbf{s}) \leq 0 \}$$

where the state of the system should always be.

- In the “MDP theory”, assign an infinite penalty to leaving \mathbb{F} , i.e. add:

$$I_{\mathbb{F}}(\mathbf{s}, \mathbf{a}) = \begin{cases} 0 & \text{if } \mathbf{s} \in \mathbb{F} \\ +\infty & \text{if } \mathbf{s} \notin \mathbb{F} \end{cases}$$

to stage cost L .

- In RL, ∞ penalties are not meaningful: “There is no backup from death”
- Common approach: assign a “very large” penalty to $\mathbf{s} \notin \mathbb{F}$ instead of $+\infty$.
- Use of “barrier functions” in RL

Why discounting?

MDP:

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

where $\mathbf{a}_k = \pi(\mathbf{s}_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

Discounting is (in general) needed to make the MDP well defined, is that all?

Can we give an interpretation of discounting?

Why discounting?

MDP:

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

where $\mathbf{a}_k = \pi(\mathbf{s}_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

Discounting is (in general) needed to make the MDP well defined, is that all?

Can we give an interpretation of discounting?

System lifetime: assuming that the system can (irremediably) fail at any time k with probability $1 - \gamma$, then discounting accounts for resulting probabilistic lifetime.

Why discounting?

MDP:

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

where $\mathbf{a}_k = \pi(s_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

Discounting is (in general) needed to make the MDP well defined, is that all?

Can we give an interpretation of discounting?

System lifetime: assuming that the system can (irremediably) fail at any time k with probability $1 - \gamma$, then discounting accounts for resulting probabilistic lifetime.

E.g. a system with a sampling time of 1 second, and a 90% chance of having a lifetime of 20 years, should have $\gamma = 0.999999996349275$

Why discounting?

MDP:

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

where $\mathbf{a}_k = \pi(s_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

Discounting is (in general) needed to make the MDP well defined, is that all?

Can we give an interpretation of discounting?

Investment model: expected economic growth r (per time unit) implies that earning at time k is worth $(1+r)^{-k}$ the same earning at time 0. Hence $\gamma = (1+r)^{-1}$.

Why discounting?

MDP:

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

where $\mathbf{a}_k = \pi(s_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

Discounting is (in general) needed to make the MDP well defined, is that all?

Can we give an interpretation of discounting?

Investment model: expected economic growth r (per time unit) implies that earning at time k is worth $(1+r)^{-k}$ the same earning at time 0. Hence $\gamma = (1+r)^{-1}$.

E.g. a system with a sampling time of 1 second and an expected return of 10% per year should have $\gamma = 0.99999999848887$

Why discounting?

MDP:

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

where $\mathbf{a}_k = \pi(s_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

Discounting is (in general) needed to make the MDP well defined, is that all?

Can we give an interpretation of discounting?

Investment model: expected economic growth r (per time unit) implies that earning at time k is worth $(1+r)^{-k}$ the same earning at time 0. Hence $\gamma = (1+r)^{-1}$.

E.g. a system with a sampling time of 1 second and an expected return of 10% per year should have $\gamma = 0.999999999848887$

Bottom line: on “engineering applications”, the discount tends to (should) be extremely close to 1

Why discounting?

Gain optimal MDP:

$$\min_{\pi} \lim_{N \rightarrow \infty} \mathbb{E}_{\pi} \left[\sum_{k=0}^N \frac{1}{N} L(s_k, \mathbf{a}_k) \right]$$

where $\mathbf{a}_k = \pi(s_k)$ and system dynamics

$$s_{k+1} \sim \mathbb{P}[\cdot | s_k, \mathbf{a}_k]$$

What about considering average cost?

Policy π

- is said to achieve “gain optimality”
- transients are irrelevant as they have no contribution in the average return
- tends to yield “bang-bang” actions until optimal steady state is reached
- is not unique!

Why discounting?

Gain optimal MDP:

$$\min_{\pi} \lim_{N \rightarrow \infty} \mathbb{E}_{\pi} \left[\sum_{k=0}^N \frac{1}{N} L(s_k, \mathbf{a}_k) \right]$$

where $\mathbf{a}_k = \pi(s_k)$ and system dynamics

$$s_{k+1} \sim \mathbb{P}[\cdot | s_k, \mathbf{a}_k]$$

What about considering average cost?

Policy π

- is said to achieve “gain optimality”
- transients are irrelevant as they have no contribution in the average return
- tends to yield “bang-bang” actions until optimal steady state is reached
- is not unique!

... gain optimal policies are of questionable use for control

Why discounting?

Bias optimal MDP:

$$\min_{\pi} \mathbb{E}_{\pi} \left[\sum_{k=0}^N L(\mathbf{s}_k, \mathbf{a}_k) - V_G^*(\mathbf{s}_0) \right]$$

where $\mathbf{a}_k = \pi(\mathbf{s}_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

What about “removing” the average cost?

where V_G^* is the value function associated to gain optimal problem.

Policy π

- is said to achieve “bias optimality”
- “best transient to gain-optimal state”
- there are RL algorithms for bias optimality

Why discounting?

Bias optimal MDP:

$$\min_{\pi} \mathbb{E}_{\pi} \left[\sum_{k=0}^N L(\mathbf{s}_k, \mathbf{a}_k) - V_G^*(\mathbf{s}_0) \right]$$

where $\mathbf{a}_k = \pi(\mathbf{s}_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

What about “removing” the average cost?

where V_G^* is the value function associated to gain optimal problem.

Policy π

- is said to achieve “bias optimality”
- “best transient to gain-optimal state”
- there are RL algorithms for bias optimality

The ideas we discuss here work for all cases. Discounted problems tend to yield “more meaningful” behavior. Discounting create some challenges for stability theory though. More on this in a bit.

Outline

- 1 The Basics
- 2 More background
- 3 Let's take a deeper dive**
- 4 Parametrization & Role of the model
- 5 RL over MPC

How does MPC model an MDP?

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

How does MPC model an MDP?

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

$$V^{\text{MPC}}(\mathbf{s}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

yields $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^*$ as by-product

How does MPC model an MDP?

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

How does MPC model an MDP?

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

MPC is **consistent**, i.e.

$$V^{\text{MPC}}(\mathbf{s}) = \min_{\mathbf{a}} Q^{\text{MPC}}(\mathbf{s}, \mathbf{a})$$

$$\pi^{\text{MPC}}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q^{\text{MPC}}(\mathbf{s}, \mathbf{a})$$

→ “**sound representation**” of MDP

How does MPC model an MDP?

Markov Decision Process:

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

MPC is optimal if:

$$\pi^{\text{MPC}}(\mathbf{s}) = \pi_{\star}(\mathbf{s})$$

for all \mathbf{s}

MPC

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

MPC is **consistent**, i.e.

$$V^{\text{MPC}}(\mathbf{s}) = \min_{\mathbf{a}} Q^{\text{MPC}}(\mathbf{s}, \mathbf{a})$$

$$\pi^{\text{MPC}}(\mathbf{s}) = \arg \min_{\mathbf{a}} Q^{\text{MPC}}(\mathbf{s}, \mathbf{a})$$

→ “sound representation” of MDP

How does MPC model an MDP?

Markov Decision Process:

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

MPC is optimal if:

$$\pi^{\text{MPC}}(s) = \pi_{\star}(s)$$

for all s

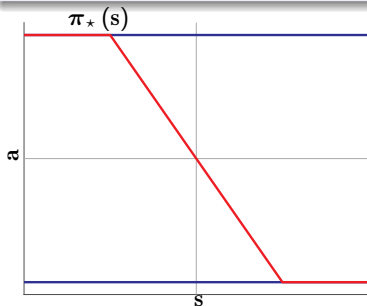
MPC

$$Q^{\text{MPC}}(s, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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



How does MPC model an MDP?

Markov Decision Process:

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

MPC is optimal if:

$$\pi^{\text{MPC}}(s) = \pi_{\star}(s)$$

for all s

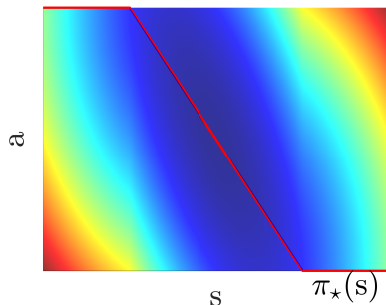
MPC

$$Q^{\text{MPC}}(s, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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



How does MPC model an MDP?

Markov Decision Process:

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

MPC is optimal if:

$$\pi^{\text{MPC}}(\mathbf{s}) = \pi_{\star}(\mathbf{s})$$

for all \mathbf{s}

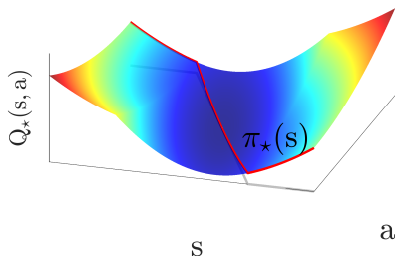
MPC

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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



How does MPC model an MDP?

Markov Decision Process:

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

MPC is optimal if:

$$\pi^{\text{MPC}}(\mathbf{s}) = \pi_{\star}(\mathbf{s})$$

for all \mathbf{s}

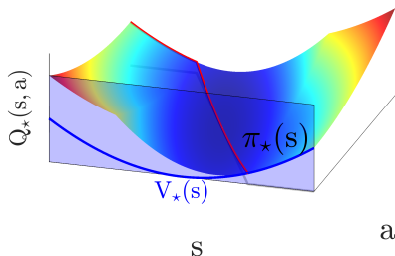
MPC

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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



How does MPC model an MDP?

Markov Decision Process:

$$\min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, a_k) \mid a_k = \pi(s_k) \right]$$

MPC is optimal if:

$$\pi^{\text{MPC}}(s) = \pi_{\star}(s)$$

for all s

MPC

$$Q^{\text{MPC}}(s, a) = \min_{x, u} T(x_N) + \sum_{k=0}^{N-1} L(x_k, u_k)$$

$$\text{s.t. } x_{k+1} = f(x_k, u_k)$$

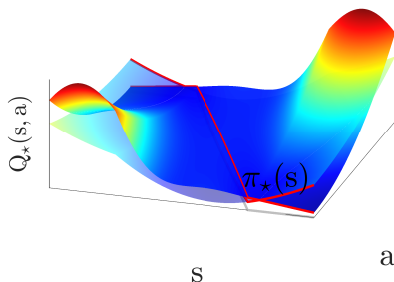
$$h(x_k, u_k) \leq 0$$

$$x_0 = s, \quad u_0 = a$$

But optimality implies only

$$\arg \max_a Q^{\text{MPC}}(s, a) = \arg \max_a Q_{\star}(s, a)$$

Optimal MPC can still be an “incomplete” model of the MDP, i.e. not a model of the value of states and actions.



How does MPC model an MDP?

Markov Decision Process:

$$\min_{\pi} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(s_k, a_k) \mid a_k = \pi(s_k) \right]$$

MPC is optimal if:

$$\pi^{\text{MPC}}(s) = \pi_{\star}(s)$$

for all s

MPC

$$Q^{\text{MPC}}(s, a) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

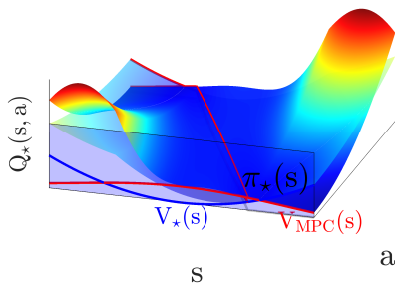
$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

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

But optimality implies only

$$\arg \max_a Q^{\text{MPC}}(s, a) = \arg \max_a Q_{\star}(s, a)$$

Optimal MPC can still be an “incomplete” model of the MDP, i.e. not a model of the value of states and actions.



How does MPC model an MDP?

Markov Decision Process:

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

MPC is a **complete** MDP model if:

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = Q_{\star}(\mathbf{s}, \mathbf{a})$$

for all \mathbf{s}, \mathbf{a}

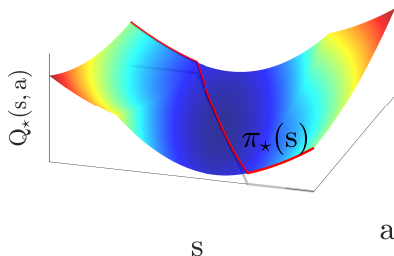
MPC

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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



How does MPC model an MDP?

Markov Decision Process:

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

MPC is a **complete** MDP model if:

$$Q^{\text{MPC}}(s, \mathbf{a}) = Q_{\star}(s, \mathbf{a})$$

for all s, \mathbf{a}

Completeness implies optimality i.e.

$$\pi^{\text{MPC}}(s) = \pi_{\star}(s)$$

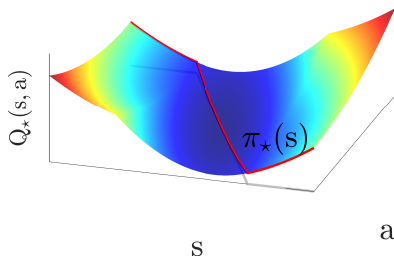
MPC

$$Q^{\text{MPC}}(s, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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



How does MPC model an MDP?

Markov Decision Process:

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

MPC is a **complete** MDP model if:

$$Q^{\text{MPC}}(s, \mathbf{a}) = Q_{\star}(s, \mathbf{a})$$

for all s, \mathbf{a}

Completeness implies optimality i.e.

$$\pi^{\text{MPC}}(s) = \pi_{\star}(s)$$

Matching the MPC action-value function to the optimal one is desirable

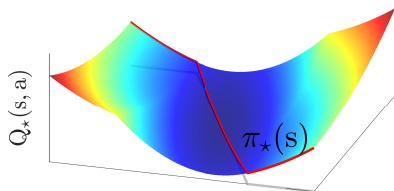
MPC

$$Q^{\text{MPC}}(s, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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



a

s

How does MPC model an MDP?

Markov Decision Process:

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

Find θ such that

$$Q_{\theta}^{\text{MPC}}(s, \mathbf{a}) = Q_{\star}(s, \mathbf{a})$$

for all s, \mathbf{a} ?

MPC

$$Q_{\theta}^{\text{MPC}}(s, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

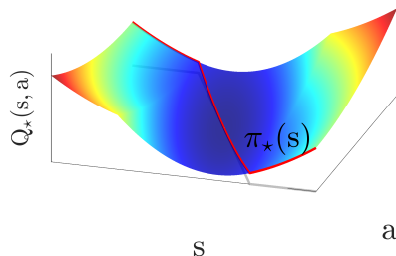
$$\text{s.t. } \mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

High demand in the MPC model !

But more on that in Lecture 3...



More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

$$V^{\text{MPC}}(\mathbf{s}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

yields $\pi^{\text{MPC}}(\mathbf{s}) = \mathbf{u}_0^{\star}$ as by-product

More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

$$Q^{\text{MPC}}(\mathbf{s}, \mathbf{a}) = \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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

In general

$$\pi^{\text{MPC}} \neq \pi_{\star}, \quad V^{\text{MPC}} \neq V_{\star}, \quad Q^{\text{MPC}} \neq Q_{\star}$$

but...

More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

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

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

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

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

In general

$$\pi^{\text{MPC}} \neq \pi_{\star}, \quad V^{\text{MPC}} \neq V_{\star}, \quad Q^{\text{MPC}} \neq Q_{\star}$$

but...

More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

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

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

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

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

Theorem: under some assumptions

$$\pi_{\theta} = \pi_{\star}, \quad V_{\theta} = V_{\star}, \quad Q_{\theta} = Q_{\star}$$

hold for some $T_{\theta}, L_{\theta}, \mathbf{h}_{\theta}$

More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

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

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

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

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

Theorem: under some assumptions

$$\pi_{\theta} = \pi_{\star}, \quad V_{\theta} = V_{\star}, \quad Q_{\theta} = Q_{\star}$$

hold for some $T_{\theta}, L_{\theta}, \mathbf{h}_{\theta}$

- MPC can “capture” $\pi_{\star}, Q_{\star}, V_{\star}$, even if MPC model is inaccurate
- Requires modifications of the stage cost & constraints
- Valid for all MPC schemes (classic, robust, stochastic, economic, etc)

More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

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

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

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

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

Theorem: under some assumptions

$$\pi_{\theta} = \pi_{\star}, \quad V_{\theta} = V_{\star}, \quad Q_{\theta} = Q_{\star}$$

hold for some $T_{\theta}, L_{\theta}, \mathbf{h}_{\theta}$

Learning+MPC where cost & constraints are adjusted is formally justified

“Holistic” view of MPC: model for Q_{\star} , cost & constraints are part of that

More on the full MPC parametrization

Markov Decision Process:

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

Value functions:

$$V_{\star}(\mathbf{s}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \right]$$

$$Q_{\star}(\mathbf{s}, \mathbf{a}) = \mathbb{E}_{\pi_{\star}} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k) \mid \mathbf{a}_0 = \mathbf{a} \right]$$

$$\pi_{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\star}(\mathbf{s}, \mathbf{a})$$

MPC

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

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

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

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

Theorem: under some assumptions

$$\pi_{\theta} = \pi_{\star}, \quad V_{\theta} = V_{\star}, \quad Q_{\theta} = Q_{\star}$$

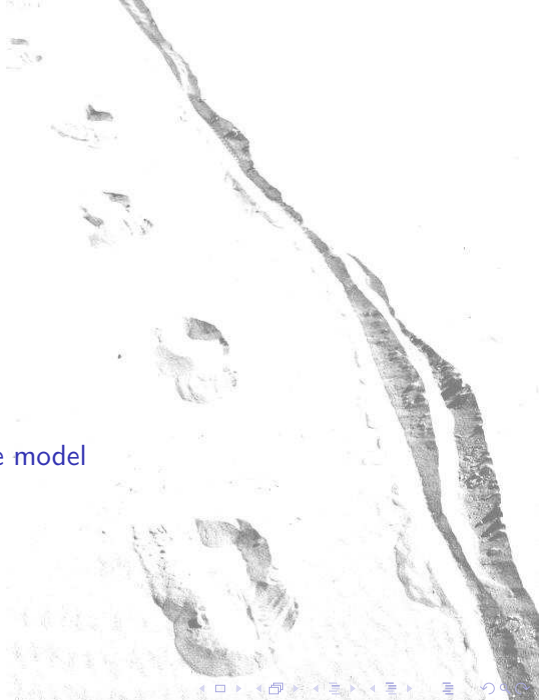
hold for some $T_{\theta}, L_{\theta}, \mathbf{h}_{\theta}$

Learning+MPC where cost & constraints are adjusted is formally justified

“Holistic” view of MPC: model for Q_{\star} , cost & constraints are part of that

Outline

- 1 The Basics
- 2 More background
- 3 Let's take a deeper dive
- 4 Parametrization & Role of the model**
- 5 RL over MPC



What MPC parametrization?

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

What MPC parametrization?

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

Theory says:

$$L_{\theta}(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \Delta(\mathbf{x}, \mathbf{u})$$

$$\Delta(\mathbf{x}, \mathbf{u}) = \underbrace{\mathbb{E}[V_{\star}(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]}_{\text{Real system}} - \underbrace{V_{\star}(\mathbf{f}_{\theta}(\mathbf{x}, \mathbf{u}))}_{\text{Model}}$$

$\mathbf{h}_{\theta} > 0 \leftrightarrow \infty$ values in modification

What MPC parametrization?

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

Theory says:

$$L_{\theta}(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \Delta(\mathbf{x}, \mathbf{u})$$

$$\Delta(\mathbf{x}, \mathbf{u}) = \underbrace{\mathbb{E}[V_{\star}(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]}_{\text{Real system}} - \underbrace{V_{\star}(\mathbf{f}_{\theta}(\mathbf{x}, \mathbf{u}))}_{\text{Model}}$$

$\mathbf{h}_{\theta} > 0 \leftrightarrow \infty$ values in modification

Remarks:

- In practice, Δ (or L_{θ}) parametrized in a chosen class of functions, and “learned”
- L_{θ} , \mathbf{h}_{θ} convex is very beneficial, maybe restrictive
- When is the model optimal as is? We will come back to that later...

What MPC parametrization?

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

Theory is not very restrictive on admissible models \mathbf{f}_{θ} . Should we be worried? Not necessarily... this theory is not the end of the story

Theory says:

$$L_{\theta}(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \Delta(\mathbf{x}, \mathbf{u})$$

$$\Delta(\mathbf{x}, \mathbf{u}) = \underbrace{\mathbb{E}[V_{\star}(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]}_{\text{Real system}} - \underbrace{V_{\star}(\mathbf{f}_{\theta}(\mathbf{x}, \mathbf{u}))}_{\text{Model}}$$

$\mathbf{h}_{\theta} > 0 \leftrightarrow \infty$ values in modification

Remarks:

- In practice, Δ (or L_{θ}) parametrized in a chosen class of functions, and “learned”
- L_{θ} , \mathbf{h}_{θ} convex is very beneficial, maybe restrictive
- When is the model optimal as is? We will come back to that later...

Role of the MPC model?

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

Role of the MPC model?

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

Remarks:

- Theory not very restrictive on admissible models \mathbf{f}_{θ}
- Examples where \mathbf{f}_{θ} is “very wrong” but MPC gives Q^*, V^*, π^*
- θ such that MPC gives Q^*, V^*, π^* may be non-unique
- Best “SYSID model” is not necessarily the best MPC model

What is the role of the model?

Role of the MPC model?

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

Remarks:

- Theory not very restrictive on admissible models \mathbf{f}_{θ}
- Examples where \mathbf{f}_{θ} is “very wrong” but MPC gives Q^*, V^*, π^*
- θ such that MPC gives Q^*, V^*, π^* may be non-unique
- Best “SYSID model” is not necessarily the best MPC model

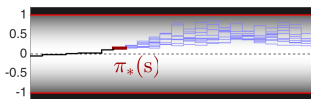
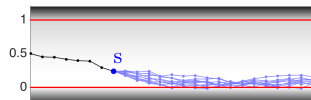
What is the role of the model?

Reflections: *depending on how we view this, it can become “philosophical”*

- MPC plan provides **explainability**. Wrong model \Rightarrow wrong plan \Rightarrow no explainability.
- MPC model associated to **safety** (*more on that soon*)

Role of the model - Explainable RL

Benefit of MPC over “black-box” RL



Role of the model - Explainable RL

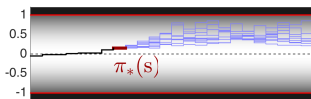
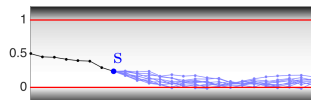
Benefit of MPC over “black-box” RL

MPC provides explainability...

... if model “makes sense”

- Not required by theory
- Not necessarily done by RL for MPC

How to keep the model sensible?



Role of the model - Explainable RL

Benefit of MPC over “black-box” RL

MPC provides explainability...

... if model “makes sense”

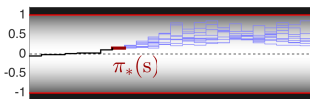
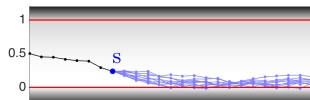
- Not required by theory
- Not necessarily done by RL for MPC

Theory says:

$$L_{\theta}(x, u) = L(x, u) + \Delta(x, u)$$

$$\Delta(x, u) = \underbrace{\mathbb{E}[V_{*}(x_{+}) \mid x, u]}_{\text{Real system}} - \underbrace{V_{*}(f_{\theta}(x, u))}_{\text{Model}}$$

How to keep the model sensible?



Role of the model - Explainable RL

Benefit of MPC over “black-box” RL

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

Theory says:

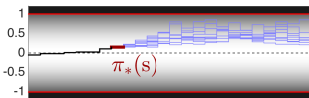
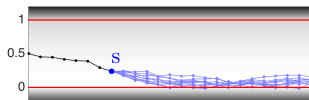
$$L_{\theta}(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \Delta(\mathbf{x}, \mathbf{u})$$

$$\Delta(\mathbf{x}, \mathbf{u}) = \underbrace{\mathbb{E}[V_{\star}(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]}_{\text{Real system}} - \underbrace{V_{\star}(\mathbf{f}_{\theta}(\mathbf{x}, \mathbf{u}))}_{\text{Model}}$$

Adjust θ such that

- $\mathbb{P}[\mathbf{f}_{\theta}(\mathbf{s}, \mathbf{a}) | \mathbf{s}, \mathbf{a}]$ (likelihood) is “high”
- $L_{\theta}, \mathbf{h}_{\theta}$ gives optimal MPC

for all \mathbf{s}, \mathbf{a} in data



Role of the model - Explainable RL

Benefit of MPC over “black-box” RL

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

Theory says:

$$L_{\theta}(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \Delta(\mathbf{x}, \mathbf{u})$$

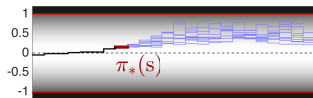
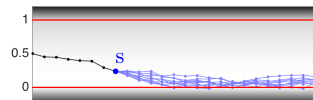
$$\Delta(\mathbf{x}, \mathbf{u}) = \underbrace{\mathbb{E}[V_{\star}(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]}_{\text{Real system}} - \underbrace{V_{\star}(\mathbf{f}_{\theta}(\mathbf{x}, \mathbf{u}))}_{\text{Model}}$$

Adjust θ such that

- $\mathbb{P}[\mathbf{f}_{\theta}(\mathbf{s}, \mathbf{a}) | \mathbf{s}, \mathbf{a}]$ (likelihood) is “high”
- $L_{\theta}, \mathbf{h}_{\theta}$ gives optimal MPC

for all \mathbf{s}, \mathbf{a} in data

RL & SYSID ought to combine without contradiction. If performance is key, RL should superseded SYSID.



Role of the model - Explainable RL

Benefit of MPC over “black-box” RL

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

Theory says:

$$L_{\theta}(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \Delta(\mathbf{x}, \mathbf{u})$$

$$\Delta(\mathbf{x}, \mathbf{u}) = \underbrace{\mathbb{E}[V_{\star}(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]}_{\text{Real system}} - \underbrace{V_{\star}(\mathbf{f}_{\theta}(\mathbf{x}, \mathbf{u}))}_{\text{Model}}$$

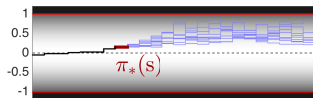
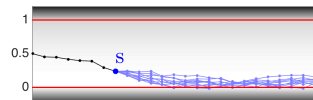
Adjust θ such that

- $\mathbb{P}[\mathbf{f}_{\theta}(\mathbf{s}, \mathbf{a}) | \mathbf{s}, \mathbf{a}]$ (likelihood) is “high”
- $L_{\theta}, \mathbf{h}_{\theta}$ gives optimal MPC

for all \mathbf{s}, \mathbf{a} in data

SYSID & RL can do

- Need to “harmonize” the two methods
- Take “SYSID steps” in the null space of $\nabla_{\theta}^2 J(\pi_{\theta}^{\text{MPC}})$



Role of the model - Explainable RL

Benefit of MPC over “black-box” RL

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

Theory says:

$$L_{\theta}(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \Delta(\mathbf{x}, \mathbf{u})$$

$$\Delta(\mathbf{x}, \mathbf{u}) = \underbrace{\mathbb{E}[V_{\star}(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]}_{\text{Real system}} - \underbrace{V_{\star}(\mathbf{f}_{\theta}(\mathbf{x}, \mathbf{u}))}_{\text{Model}}$$

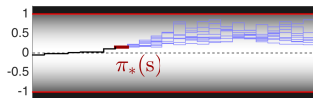
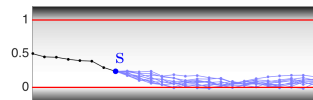
Adjust θ such that

- $\mathbb{P}[\mathbf{f}_{\theta}(\mathbf{s}, \mathbf{a}) | \mathbf{s}, \mathbf{a}]$ (likelihood) is “high”
- $L_{\theta}, \mathbf{h}_{\theta}$ gives optimal MPC

for all \mathbf{s}, \mathbf{a} in data

Reflection:

Do we need a concept of “explainability” for MPC?
What fundamental properties should the MPC model have to be deemed “explaining”?



Outline

- 1 The Basics
- 2 More background
- 3 Let's take a deeper dive
- 4 Parametrization & Role of the model
- 5 RL over MPC

Classic RL vs. RL-MPC

MDP:

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

where $\mathbf{a}_k = \pi(\mathbf{s}_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

MPC:

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

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

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

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

yields π^{MPC} , V^{MPC} , and Q^{MPC}

Classic RL vs. RL-MPC

MDP:

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

where $\mathbf{a}_k = \pi(\mathbf{s}_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

MPC:

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

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

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

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

yields π^{MPC} , V^{MPC} , and Q^{MPC}

RL with DNN

- correct structure is unknown
- good initialization is difficult
- respecting constraints is difficult & implicit

Classic RL vs. RL-MPC

MDP:

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

where $\mathbf{a}_k = \pi(\mathbf{s}_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

MPC:

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

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

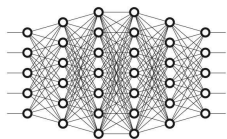
$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

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

yields π^{MPC} , V^{MPC} , and Q^{MPC}

RL with DNN

- correct structure is unknown
- good initialization is difficult
- respecting constraints is difficult & implicit



Classic RL vs. RL-MPC

MDP:

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

where $\mathbf{a}_k = \pi(\mathbf{s}_k)$ and system dynamics

$$\mathbf{s}_{k+1} \sim \mathbb{P}[\cdot | \mathbf{s}_k, \mathbf{a}_k]$$

MPC:

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

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

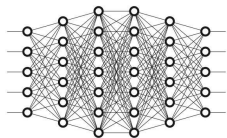
$$\mathbf{h}_{\theta}(\mathbf{x}_k, \mathbf{u}_k) \leq 0$$

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

yields π^{MPC} , V^{MPC} , and Q^{MPC}

RL with DNN

- correct structure is unknown
- good initialization is difficult
- respecting constraints is difficult & implicit



MPC

- Structure and initialization given
- Constraints enforced explicitly
- Theory says that we can get V_* , Q_* , π_* from MPC

RL methods & MPC

Form function approximators:

$$Q_{\theta}(s, a), V_{\theta}(s), \pi_{\theta}(s)$$

via ad-hoc parametrization

Form function approximators:

$$Q_{\theta}(s, \mathbf{a}), V_{\theta}(s), \pi_{\theta}(s)$$

via ad-hoc parametrization

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

$$Q_{\theta}(s, \mathbf{a}) \approx Q_{*}(s, \mathbf{a})$$

Yields policy:

$$\pi_{\theta}(s) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\theta}(s, \mathbf{a}) \approx \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{*}(s, \mathbf{a}) = \pi_{*}(s)$$

E.g. basic Q-learning uses:

$$\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} Q_{\theta}(s_k, \mathbf{a}_k)$$

$$\delta = L(s_k, \mathbf{a}_k) + \gamma V_{\theta}(s_{k+1}) - Q_{\theta}(s_k, \mathbf{a}_k)$$

Form function approximators:

$$Q_{\theta}(s, \mathbf{a}), V_{\theta}(s), \pi_{\theta}(s)$$

via ad-hoc parametrization

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

$$Q_{\theta}(s, \mathbf{a}) \approx Q_{*}(s, \mathbf{a})$$

Yields policy:

$$\pi_{\theta}(s) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\theta}(s, \mathbf{a}) \approx \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{*}(s, \mathbf{a}) = \pi_{*}(s)$$

E.g. basic Q-learning uses:

$$\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} Q_{\theta}(s_k, \mathbf{a}_k)$$

$$\delta = L(s_k, \mathbf{a}_k) + \gamma V_{\theta}(s_{k+1}) - Q_{\theta}(s_k, \mathbf{a}_k)$$

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

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

yields policy $\pi_{\theta}(\mathbf{x}) \approx \pi_{*}(\mathbf{x})$ directly.

Form function approximators:

$$Q_{\theta}(s, \mathbf{a}), V_{\theta}(s), \pi_{\theta}(s)$$

via ad-hoc parametrization

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

$$Q_{\theta}(s, \mathbf{a}) \approx Q_{*}(s, \mathbf{a})$$

Yields policy:

$$\pi_{\theta}(s) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\theta}(s, \mathbf{a}) \approx \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{*}(s, \mathbf{a}) = \pi_{*}(s)$$

E.g. basic Q-learning uses:

$$\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} Q_{\theta}(s_k, \mathbf{a}_k)$$

$$\delta = L(s_k, \mathbf{a}_k) + \gamma V_{\theta}(s_{k+1}) - Q_{\theta}(s_k, \mathbf{a}_k)$$

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

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

yields policy $\pi_{\theta}(\mathbf{x}) \approx \pi_{*}(\mathbf{x})$ directly. E.g.

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}[\nabla_{\theta} \pi_{\theta} \nabla_{\mathbf{a}} Q_{\pi_{\theta}}]$$

- **Derivative-free methods**

- ▶ Build a surrogate of $J(\pi_{\theta})$
- ▶ Optimize over that model
- ▶ Difficult over large parameter spaces

RL methods & MPC

Form function approximators:

$$Q_{\theta}(s, a), V_{\theta}(s), \pi_{\theta}(s)$$

via ad-hoc parametrization

Derivative-based methods
require Q_{θ} , V_{θ} , π_{θ} and
computing their **sensitivities**

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

$$Q_{\theta}(s, a) \approx Q_{*}(s, a)$$

Yields policy:

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

E.g. basic Q-learning uses:

$$\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} Q_{\theta}(s_k, a_k)$$

$$\delta = L(s_k, a_k) + \gamma V_{\theta}(s_{k+1}) - Q_{\theta}(s_k, a_k)$$

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

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

yields policy $\pi_{\theta}(x) \approx \pi_{*}(x)$ directly. E.g.

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}[\nabla_{\theta} \pi_{\theta} \nabla_a Q_{\pi_{\theta}}]$$

- **Derivative-free methods**

- ▶ Build a surrogate of $J(\pi_{\theta})$
- ▶ Optimize over that model
- ▶ Difficult over large parameter spaces

RL methods & MPC

Form function approximators:

$$Q_{\theta}(s, \mathbf{a}), V_{\theta}(s), \pi_{\theta}(s)$$

via ad-hoc parametrization

Derivative-based methods require Q_{θ} , V_{θ} , π_{θ} and computing their **sensitivities**

In the RL-MPC context, Q_{θ} , V_{θ} , π_{θ} are coming from an MPC scheme, typically cast as **Nonlinear Program**. What about the sensitivities?

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

$$Q_{\theta}(s, \mathbf{a}) \approx Q_{*}(s, \mathbf{a})$$

Yields policy:

$$\pi_{\theta}(s) = \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{\theta}(s, \mathbf{a}) \approx \underset{\mathbf{a}}{\operatorname{arg\,min}} Q_{*}(s, \mathbf{a}) = \pi_{*}(s)$$

E.g. basic Q-learning uses:

$$\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} Q_{\theta}(s_k, \mathbf{a}_k)$$

$$\delta = L(s_k, \mathbf{a}_k) + \gamma V_{\theta}(s_{k+1}) - Q_{\theta}(s_k, \mathbf{a}_k)$$

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

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

yields policy $\pi_{\theta}(\mathbf{x}) \approx \pi_{*}(\mathbf{x})$ directly. E.g.

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}[\nabla_{\theta} \pi_{\theta} \nabla_{\mathbf{a}} Q_{\pi_{\theta}}]$$

- **Derivative-free methods**

- ▶ Build a surrogate of $J(\pi_{\theta})$
- ▶ Optimize over that model
- ▶ Difficult over large parameter spaces

MPC Sensitivities?

MPC is a Nonlinear Program

Optimal value

$$\begin{aligned} V_{\theta}(s) &= \min_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad &\mathbf{g}(\mathbf{w}, s, \theta) = 0 \\ &\mathbf{h}(\mathbf{w}, s, \theta) \leq 0 \end{aligned}$$

Optimal solution

$$\begin{aligned} \mathbf{w}_{\theta}^*(s) &= \text{a min}_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad &\dots \end{aligned}$$

MPC Sensitivities?

MPC is a Nonlinear Program

Optimal value

$$\begin{aligned} V_{\theta}(s) &= \min_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t. } \mathbf{g}(\mathbf{w}, s, \theta) &= 0 \\ \mathbf{h}(\mathbf{w}, s, \theta) &\leq 0 \end{aligned}$$

Optimal solution

$$\begin{aligned} \mathbf{w}_{\theta}^*(s) &= \underset{\mathbf{w}}{\text{a min}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t. } &\dots \end{aligned}$$

How to obtain:

$$\nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}, \nabla_{\theta} \mathbf{w}_{\theta}^*$$

MPC Sensitivities?

MPC is a Nonlinear Program

Optimal value

$$\begin{aligned} V_{\theta}(s) &= \min_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t. } \mathbf{g}(\mathbf{w}, s, \theta) &= 0 \\ \mathbf{h}(\mathbf{w}, s, \theta) &\leq 0 \end{aligned}$$

Optimal solution

$$\begin{aligned} \mathbf{w}_{\theta}^*(s) &= \text{a min}_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t. } &\dots \end{aligned}$$

NLP solution satisfies (KKT conditions)

$$\mathbf{r} = \begin{bmatrix} \nabla_{\mathbf{w}} \mathcal{L} \\ \mathbf{g} \\ \mathbf{h}_i \mu_i \end{bmatrix} = 0$$
$$\mathbf{h} \leq 0, \mu \geq 0$$

where Lagrange function is

$$\mathcal{L} = \Phi + \boldsymbol{\lambda}^{\top} \mathbf{g} + \boldsymbol{\mu}^{\top} \mathbf{h}$$

and $\boldsymbol{\lambda}, \boldsymbol{\mu}$ are the dual variables

How to obtain:

$$\nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}, \nabla_{\theta} \mathbf{w}_{\theta}^*$$

MPC Sensitivities?

MPC is a Nonlinear Program

Optimal value

$$\begin{aligned} V_{\theta}(s) &= \min_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad &\mathbf{g}(\mathbf{w}, s, \theta) = 0 \\ &\mathbf{h}(\mathbf{w}, s, \theta) \leq 0 \end{aligned}$$

Optimal solution

$$\begin{aligned} \mathbf{w}_{\theta}^*(s) &= \text{a min}_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad &\dots \end{aligned}$$

How to obtain:

$$\nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}, \nabla_{\theta} \mathbf{w}_{\theta}^*$$

NLP solution satisfies (KKT conditions)

$$\mathbf{r} = \begin{bmatrix} \nabla_{\mathbf{w}} \mathcal{L} \\ \mathbf{g} \\ \mathbf{h}_i \mu_i \end{bmatrix} = 0$$
$$\mathbf{h} \leq 0, \mu \geq 0$$

where Lagrange function is

$$\mathcal{L} = \Phi + \boldsymbol{\lambda}^{\top} \mathbf{g} + \boldsymbol{\mu}^{\top} \mathbf{h}$$

and $\boldsymbol{\lambda}, \boldsymbol{\mu}$ are the dual variables

MPC Sensitivities?

MPC is a Nonlinear Program

Optimal value

$$\begin{aligned} V_{\theta}(s) = \min_{\mathbf{w}} \quad & \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad & \mathbf{g}(\mathbf{w}, s, \theta) = 0 \\ & \mathbf{h}(\mathbf{w}, s, \theta) \leq 0 \end{aligned}$$

Optimal solution

$$\begin{aligned} \mathbf{w}_{\theta}^*(s) = \underset{\mathbf{w}}{\text{a min}} \quad & \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad & \dots \end{aligned}$$

How to obtain:

$$\nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}, \nabla_{\theta} \mathbf{w}_{\theta}^*$$

NLP solution satisfies (KKT conditions)

$$\mathbf{r} = \begin{bmatrix} \nabla_{\mathbf{w}} \mathcal{L} \\ \mathbf{g} \\ \mathbf{h}_i \mu_i \end{bmatrix} = 0$$
$$\mathbf{h} \leq 0, \mu \geq 0$$

where Lagrange function is

$$\mathcal{L} = \Phi + \boldsymbol{\lambda}^{\top} \mathbf{g} + \boldsymbol{\mu}^{\top} \mathbf{h}$$

and $\boldsymbol{\lambda}, \boldsymbol{\mu}$ are the dual variables

Solve NLP for s, θ , provides $\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}$, then:

$$\nabla_{\theta} V_{\theta}(s) = \nabla_{\theta} \mathcal{L}(\mathbf{w}, s, \theta, \boldsymbol{\lambda}, \boldsymbol{\mu})$$

is a simple function evaluation

MPC Sensitivities?

MPC is a Nonlinear Program

Optimal value

$$\begin{aligned} V_{\theta}(s) &= \min_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad &\mathbf{g}(\mathbf{w}, s, \theta) = 0 \\ &\mathbf{h}(\mathbf{w}, s, \theta) \leq 0 \end{aligned}$$

Optimal solution

$$\begin{aligned} \mathbf{w}_{\theta}^*(s) &= \mathop{\text{a min}}_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad &\dots \end{aligned}$$

How to obtain:

$$\nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}, \nabla_{\theta} \mathbf{w}_{\theta}^*$$

NLP solution satisfies (KKT conditions)

$$\mathbf{r} = \begin{bmatrix} \nabla_{\mathbf{w}} \mathcal{L} \\ \mathbf{g} \\ \mathbf{h}_i \mu_i \end{bmatrix} = 0$$
$$\mathbf{h} \leq 0, \mu \geq 0$$

where Lagrange function is

$$\mathcal{L} = \Phi + \boldsymbol{\lambda}^{\top} \mathbf{g} + \boldsymbol{\mu}^{\top} \mathbf{h}$$

and $\boldsymbol{\lambda}, \boldsymbol{\mu}$ are the dual variables

Solve NLP for s, θ , provides $\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}$, then:

$$\frac{\partial \mathbf{w}_{\theta}^*}{\partial \theta} = - \frac{\partial \mathbf{r}}{\partial \mathbf{w}}^{-1} \frac{\partial \mathbf{r}}{\partial \theta}$$

where $\frac{\partial \mathbf{r}}{\partial \mathbf{w}}^{-1}$ is already built in the solver, works if LICQ / SOSC

MPC Sensitivities?

MPC is a Nonlinear Program

Optimal value

$$\begin{aligned} V_{\theta}(s) = \min_{\mathbf{w}} \quad & \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad & \mathbf{g}(\mathbf{w}, s, \theta) = 0 \\ & \mathbf{h}(\mathbf{w}, s, \theta) \leq 0 \end{aligned}$$

Optimal solution

$$\begin{aligned} \mathbf{w}_{\theta}^*(s) = \underset{\mathbf{w}}{\text{a min}} \quad & \Phi(\mathbf{w}, s, \theta) \\ \text{s.t.} \quad & \dots \end{aligned}$$

How to obtain:

$$\nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}, \nabla_{\theta} \mathbf{w}_{\theta}^*$$

NLP solution satisfies (KKT conditions)

$$\mathbf{r} = \begin{bmatrix} \nabla_{\mathbf{w}} \mathcal{L} \\ \mathbf{g} \\ \mathbf{h}_i \mu_i \end{bmatrix} = 0$$
$$\mathbf{h} \leq 0, \mu \geq 0$$

where Lagrange function is

$$\mathcal{L} = \Phi + \boldsymbol{\lambda}^{\top} \mathbf{g} + \boldsymbol{\mu}^{\top} \mathbf{h}$$

and $\boldsymbol{\lambda}, \boldsymbol{\mu}$ are the dual variables

**Sensitivities do not exist for all s, \mathbf{a} .
Does that matter?**

MPC Sensitivities?

MPC is a Nonlinear Program

Optimal value

$$\begin{aligned} V_{\theta}(s) &= \min_{\mathbf{w}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t. } \mathbf{g}(\mathbf{w}, s, \theta) &= 0 \\ \mathbf{h}(\mathbf{w}, s, \theta) &\leq 0 \end{aligned}$$

Optimal solution

$$\begin{aligned} \mathbf{w}_{\theta}^*(s) &= \underset{\mathbf{w}}{\text{a min}} \Phi(\mathbf{w}, s, \theta) \\ \text{s.t. } &\dots \end{aligned}$$

How to obtain:

$$\nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}, \nabla_{\theta} \mathbf{w}_{\theta}^*$$

NLP solution satisfies (KKT conditions)

$$\mathbf{r} = \begin{bmatrix} \nabla_{\mathbf{w}} \mathcal{L} \\ \mathbf{g} \\ \mathbf{h}_i \mu_i \end{bmatrix} = 0$$
$$\mathbf{h} \leq 0, \mu \geq 0$$

where Lagrange function is

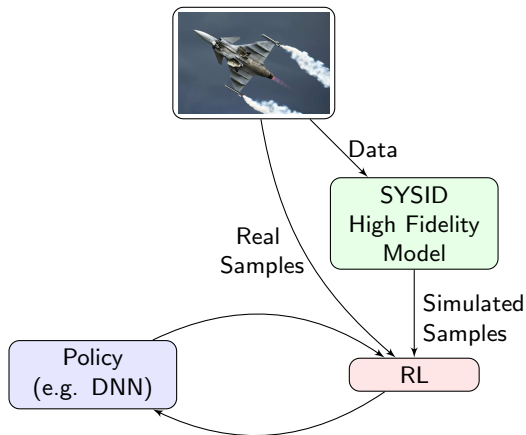
$$\mathcal{L} = \Phi + \boldsymbol{\lambda}^{\top} \mathbf{g} + \boldsymbol{\mu}^{\top} \mathbf{h}$$

and $\boldsymbol{\lambda}, \boldsymbol{\mu}$ are the dual variables

**Sensitivities do not exist for all s , a.
Does that matter?**

In general no: they exist *almost everywhere*, and always appear inside $\mathbb{E}[\cdot]$. If the MDP has underlying densities, then we are good.

Model-based RL methods vs. RL-MPC: Data flow



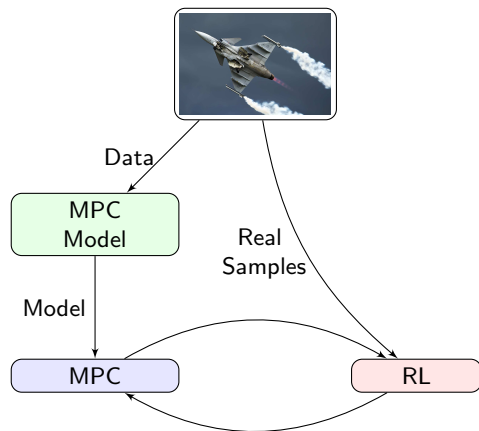
Common setup for “classic RL:

- Build statistical model of the real system
- Generate simulated samples
- Feed RL with real and simulated samples

Remarks:

- Simulated data much cheaper than real ones, most data will be simulated ones
- With mostly simulated data:
 - ▶ \approx equivalent to approximate DP
 - ▶ policy optimality relies on model quality

Model-based RL methods vs. RL-MPC: Data flow



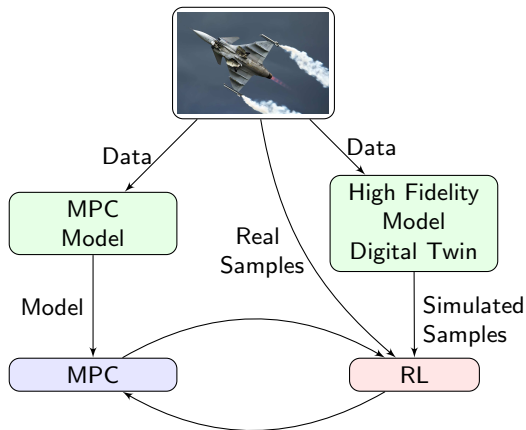
Basic setup for “RL-MPC”:

- Build MPC model of the real system
- Pass it to MPC scheme
- Feed RL with real samples

Remarks:

- RL tunes MPC for real system
- MPC model may be “detuned” from SYSID version
- Real data are expensive...

Model-based RL methods vs. RL-MPC: Data flow



“Mixed” setup for “RL-MPC”:

- Build MPC model of the real system
- MPC model is typically “simple”
- Build statistical model of the real system
- Generate simulated samples
- Feed RL with real and simulated samples

Remarks:

- Simple MPC model
- Complex simulation model
- MPC model may be “detuned” from SYSID version

- MPC as a path for safety and stability in RL
- More results & ideas

Thanks for your attention!