So far we have seen...

- The building blocks
 - MDPs, DP, Bellman, Policy and Value Iteration, Reinforcement Learning, Actor-Critic
 - LQR, Dynamic Optimization, Nonlinear Programming, MPC, Sensitivities
- MPC & RL
 - Synthesis Different flavors of ML / RL and MPC
 - Software tool leap-c

S. Gros (NTNU) Intro to RL-MPC Fall 2025 1/29

So far we have seen...

- The building blocks
 - MDPs, DP, Bellman, Policy and Value Iteration, Reinforcement Learning, Actor-Critic
 - LQR, Dynamic Optimization, Nonlinear Programming, MPC, Sensitivities
- MPC & RL
 - Synthesis Different flavors of ML / RL and MPC
 - Software tool leap-c

What we will discuss: parametrized / hierarchical RL over MPC

• The theory that supports it and what does it tell us? (Now)

1/29

So far we have seen...

- The building blocks
 - MDPs, DP, Bellman, Policy and Value Iteration, Reinforcement Learning, Actor-Critic
 - LQR, Dynamic Optimization, Nonlinear Programming, MPC, Sensitivities
- MPC & RL
 - Synthesis Different flavors of ML / RL and MPC
 - Software tool leap-c

What we will discuss: parametrized / hierarchical RL over MPC

- The theory that supports it and what does it tell us? (Now)
- Safe & stable RL over MPC (In the afternoon)

1/29

So far we have seen...

- The building blocks
 - MDPs, DP, Bellman, Policy and Value Iteration, Reinforcement Learning, Actor-Critic
 - LQR, Dynamic Optimization, Nonlinear Programming, MPC, Sensitivities
- MPC & RL
 - Synthesis Different flavors of ML / RL and MPC
 - Software tool leap-c

What we will discuss: parametrized / hierarchical RL over MPC

- The theory that supports it and what does it tell us? (Now)
- Safe & stable RL over MPC (In the afternoon)
- RL over MPC with belief states a future prospect (In the afternoon)

1/29

So far we have seen...

- The building blocks
 - MDPs, DP, Bellman, Policy and Value Iteration, Reinforcement Learning, Actor-Critic
 - LQR, Dynamic Optimization, Nonlinear Programming, MPC, Sensitivities
- MPC & RL

S. Gros (NTNU)

- Synthesis Different flavors of ML / RL and MPC
- Software tool leap-c

What we will discuss: parametrized / hierarchical RL over MPC

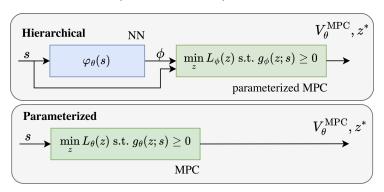
- The theory that supports it and what does it tell us? (Now)
- Safe & stable RL over MPC (In the afternoon)
- RL over MPC with belief states a future prospect (In the afternoon)
- Beyond MPC Model-based Decisions and AI for decisions (Tomorrow)

4 □ > 4 ፬ > 4 ፬ > 4 ፬ > 1 호 · 9

Fall 2025

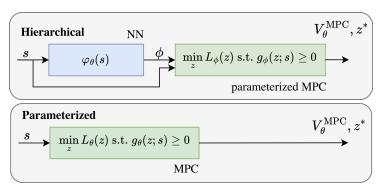
1/29

Focus of today's lectures (Leap-c structure)



2/29

Focus of today's lectures (Leap-c structure)



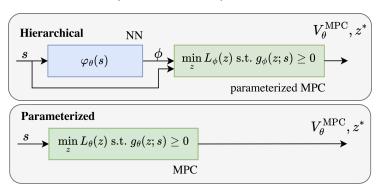
Normal thinking when using MPC

- Fit MPC model to reality as good as possible (SYSID)
- MPC cost is what we want to minimize (energy, time, money, reference)
- MPC state constraints are what we need to respect (safety)

◆ロト ◆個ト ◆意ト ◆意ト ・意 ・ からで

2 / 29

Focus of today's lectures (Leap-c structure)



Normal thinking when using MPC

- Fit MPC model to reality as good as possible (SYSID)
- MPC cost is what we want to minimize (energy, time, money, reference)
- MPC state constraints are what we need to respect (safety)

We are talking about changing everything!!

Reinforcement Learning Over MPC Why Does It Work?

Sebastien Gros

Dept. of Cybernetic, NTNU Faculty of Information Tech.

Freiburg PhD School

Outline

- 1 Let's rebuild some background MPC and MDPs
- 2 MPC as a solution to MDPs
- 3 When is RL (most) beneficial for MPC?
- 4 A Deeper Look at the Theory

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

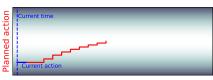
MPC: at current state s

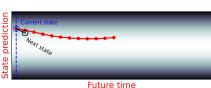
$$\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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





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

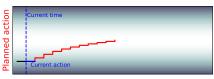
MPC: at current state s

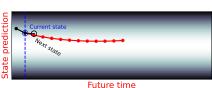
$$\min_{\mathbf{x},\mathbf{u}} T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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





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

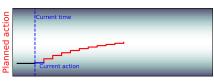
MPC: at current state s

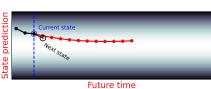
$$\min_{\mathbf{x},\mathbf{u}} \quad \mathcal{T}(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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





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

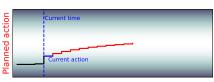
MPC: at current state s

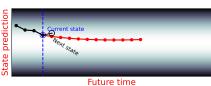
$$\min_{\mathbf{x},\mathbf{u}} \quad \mathcal{T}(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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





S. Gros (NTNU)

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

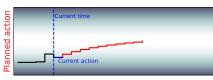
MPC: at current state s

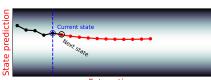
$$\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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





Future time

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

MPC: at current state s

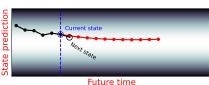
$$\min_{\mathbf{x},\mathbf{u}} \quad \mathcal{T}(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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





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

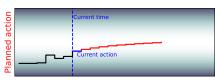
MPC: at current state s

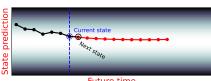
$$\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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





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

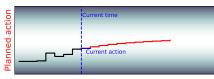
MPC: at current state s

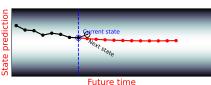
$$\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

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





S. Gros (NTNU)

Intro to RL-MPC

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

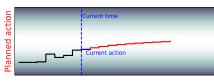
MPC: at current state s

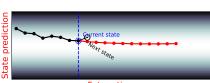
$$\min_{\mathbf{x},\mathbf{u}} \quad T\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$$

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 $\mathbf{a}=\mathbf{u}_0^{\star}$ to the system





Future time

MPC

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

$$oldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

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

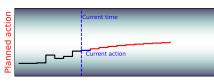
MPC: at current state s

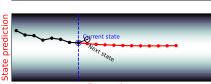
$$\min_{\mathbf{x},\mathbf{u}} \quad T\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$$

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

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

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





Future time

MPC

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

$$oldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

What an odd thing to do: we build and throw away plans all the time knowing that they are wrong all along, but kinda use them anyway...

4日 > 4目 > 4目 > 4目 > 4目 > 99の

5/29

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

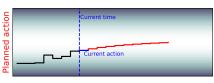
MPC: at current state s

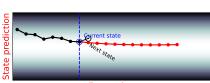
$$\min_{\mathbf{x},\mathbf{u}} \quad T\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$$

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

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

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





Future time

MPC

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

$$oldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

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

5/29

Model Predictive Control

- Model driven
- Policies from planning
- Constraints oriented

Reinforcement Learning (RL)

- Data driven
- Optimal policies from learning
- Performance oriented

6/29

Model Predictive Control

- Model driven
- Policies from planning
- Constraints oriented

Connection?

Reinforcement Learning (RL)

- Data driven
- Optimal policies from learning
- Performance oriented

6/29

Model Predictive Control

- Model driven
- Policies from planning
- Constraints oriented

Connection?

Reinforcement Learning (RL)

- Data driven
- Optimal policies from learning
- Performance oriented

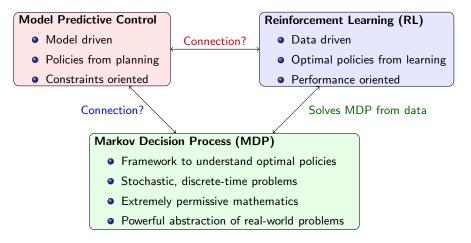
Markov Decision Process (MDP)

- Framework to understand optimal policies
- Stochastic, discrete-time problems
- Extremely permissive mathematics
- Powerful abstraction of real-world problems

Model Predictive Control Reinforcement Learning (RL) Data driven Model driven Connection? Policies from planning Optimal policies from learning Constraints oriented Performance oriented Solves MDP from data Markov Decision Process (MDP) • Framework to understand optimal policies Stochastic, discrete-time problems Extremely permissive mathematics

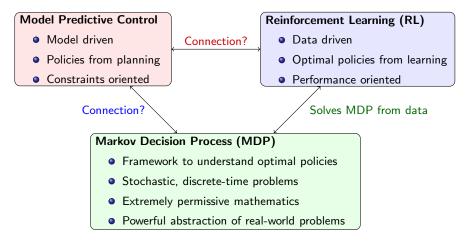
Powerful abstraction of real-world problems

Model Predictive Control Reinforcement Learning (RL) Data driven Model driven Connection? Policies from planning Optimal policies from learning Constraints oriented Performance oriented Connection? Solves MDP from data Markov Decision Process (MDP) • Framework to understand optimal policies Stochastic, discrete-time problems Extremely permissive mathematics Powerful abstraction of real-world problems



We have connected MPC and RL from an "implementation" point of view

S. Gros (NTNU) Intro to RL-MPC Fall 2025 6 / 29



We have connected MPC and RL from an "implementation" point of view But understanding what we are doing is about connecting MPC to MDPs!!

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$

S. Gros (NTNU) Intro to RL-MPC Fall 2025 7/29

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

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

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$

A (fairly) general way of describing optimal control

S. Gros (NTNU) Intro to RL-MPC Fall 2025 7 / 29

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

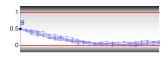
$$J(\pi) = \mathbb{E}\left[\left.\sum_{k=0}^{\infty} \gamma^k L\left(\mathbf{s}_k, \mathbf{a}_k\right) \,\right| \, \pi\right]$$

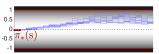
with discount $\gamma \in [0, 1]$

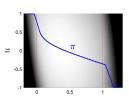
• Optimal policy: π^* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

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

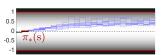
with discount $\gamma \in [0, 1]$

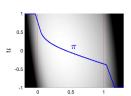
• Optimal policy: π^* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

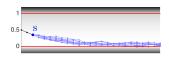
$$J(\pi) = \mathbb{E}\left[\left.\sum_{k=0}^{\infty} \gamma^k L\left(\mathbf{s}_k, \mathbf{a}_k\right) \,\right| \, \pi\right]$$

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

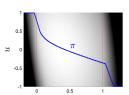
• Optimal policy: π^* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

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

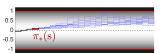
with discount $\gamma \in [0, 1]$

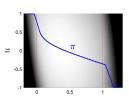
• Optimal policy: π^* from

$$\min_{x \in \mathcal{X}} J(x)$$

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

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

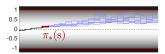
with discount $\gamma \in [0, 1]$

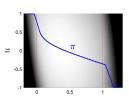
Optimal policy: π* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

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

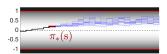
with discount $\gamma \in [0, 1]$

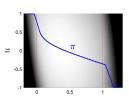
Optimal policy: π* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







7/29

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

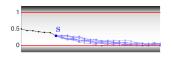
$$J(\pi) = \mathbb{E}\left[\left.\sum_{k=0}^{\infty} \gamma^k L\left(\mathbf{s}_k, \mathbf{a}_k
ight)\right| \, \pi
ight]$$

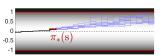
with discount $\gamma \in [0, 1]$

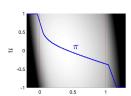
Optimal policy: π* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

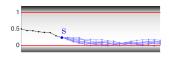
$$J(\pi) = \mathbb{E}\left[\left.\sum_{k=0}^{\infty} \gamma^k L\left(\mathbf{s}_k, \mathbf{a}_k\right) \,\right| \, \pi\right]$$

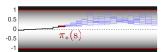
with discount $\gamma \in [0, 1]$

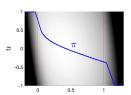
Optimal policy: π* from

$$\min_{x \in \mathcal{X}} J(x)$$

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







7/29

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

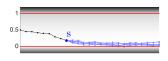
$$J(\pi) = \mathbb{E}\left[\left.\sum_{k=0}^{\infty} \gamma^k L\left(\mathbf{s}_k, \mathbf{a}_k
ight)\right| \, \pi
ight]$$

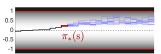
with discount $\gamma \in [0, 1]$

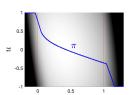
Optimal policy: π* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

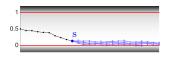
$$J(\pi) = \mathbb{E}\left[\left.\sum_{k=0}^{\infty} \gamma^k L\left(\mathbf{s}_k, \mathbf{a}_k\right) \,\right| \, \pi\right]$$

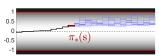
with discount $\gamma \in [0, 1]$

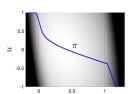
Optimal policy: π* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







7/29

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

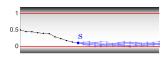
$$J(\pi) = \mathbb{E}\left[\left.\sum_{k=0}^{\infty} \gamma^k L\left(\mathbf{s}_k, \mathbf{a}_k\right) \,\right| \, \pi\right]$$

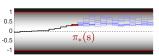
with discount $\gamma \in [0, 1]$

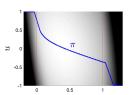
Optimal policy: π* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$







7/29

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

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

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

• Optimal policy: π^* from

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

A (fairly) general way of describing optimal control

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$

Impose hard constraints $h(s, a) \le 0$?

$$\label{eq:loss_loss} L\left(\mathbf{s},\mathbf{a}\right) = \left\{ \begin{array}{ccc} \ell\left(\mathbf{s},\mathbf{a}\right) & \text{if} & \mathbf{h}\left(\mathbf{s},\mathbf{a}\right) \leq 0 \\ \infty & \text{if} & \mathbf{h}\left(\mathbf{s},\mathbf{a}\right) > 0 \end{array} \right.$$

I will use the same view in MPC for a bit

S. Gros (NTNU) Intro to RL-MPC Fall 2025 7/29

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

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

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

• Optimal policy: π^* from

$$\min_{x \in \mathcal{X}} J(x)$$

A (fairly) general way of describing optimal control

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$

Impose hard constraints $h(s, a) \le 0$?

$$\label{eq:loss_loss} \textit{L}\left(\mathbf{s},\mathbf{a}\right) = \left\{ \begin{array}{ccc} \ell\left(\mathbf{s},\mathbf{a}\right) & \text{if} & \mathbf{h}\left(\mathbf{s},\mathbf{a}\right) \leq 0 \\ \infty & \text{if} & \mathbf{h}\left(\mathbf{s},\mathbf{a}\right) > 0 \end{array} \right.$$

I will use the same view in MPC for a bit

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

Stochastic state transitions (real world)

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

(state-action \rightarrow next state)

Policy

$$a = \pi(s)$$

is how we act on the system

Closed-loop performance

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

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

Optimal policy: π* from

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

Cost function (instant performance) $L(\mathbf{s}, \mathbf{a}) \in \mathbb{R}$

Impose hard constraints $h(s, a) \le 0$?

$$\label{eq:loss_loss} L\left(\mathbf{s},\mathbf{a}\right) = \left\{ \begin{array}{ccc} \ell\left(\mathbf{s},\mathbf{a}\right) & \text{if} & \mathbf{h}\left(\mathbf{s},\mathbf{a}\right) \leq 0 \\ \infty & \text{if} & \mathbf{h}\left(\mathbf{s},\mathbf{a}\right) > 0 \end{array} \right.$$

I will use the same view in MPC for a bit

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

S. Gros (NTNU) Intro to RL-MPC Fall 2025 7 / 29

Infinite horizon & discounted

$$oldsymbol{\pi}_{\infty}^{\star} = rg\min_{oldsymbol{\pi}} \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{a}_k)
ight]$$

Policy π : state \to action belongs to a function space

Let's transform an MDP into an MPC and understand the approximations we make

S. Gros (NTNU) Intro to RL-MPC Fall 2025 8/29

Infinite horizon & discounted

$$oldsymbol{\pi}_{\infty}^{\star} = rg \min_{oldsymbol{\pi}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{a}_k)
ight]$$

Policy π : state \rightarrow action belongs to a function space

Let's transform an MDP into an MPC and understand the approximations we make

Finite-horizon equivalent:

$$\boldsymbol{\pi}_{0,\dots,N-1}^{\star} = \operatorname*{arg\,min}_{\boldsymbol{\pi}_{0,\dots,N-1}} \mathbb{E}\left[T(\mathbf{s}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k}, \mathbf{a}_{k})\right]$$

If
$$T=V^\star$$
, then $oldsymbol{\pi}_{0,\dots,N-1}^\star=oldsymbol{\pi}_\infty^\star$

S. Gros (NTNU) Intro to RL-MPC

Infinite horizon & discounted

$$oldsymbol{\pi}_{\infty}^{\star} = rg \min_{oldsymbol{\pi}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{a}_k)
ight]$$

Policy π : state \to action belongs to a function space

Let's transform an MDP into an MPC and understand the approximations we make

Finite-horizon equivalent:

$$\boldsymbol{\pi}_{0,...,N-1}^{\star} = \operatorname*{arg\,min}_{\boldsymbol{\pi}_{0,...,N-1}} \mathbb{E}\left[T(\mathbf{s}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k}, \mathbf{a}_{k})\right]$$

If
$$T=V^\star$$
, then $oldsymbol{\pi}_{0,\dots,N-1}^\star=oldsymbol{\pi}_\infty^\star$

Planning instead of a policy:

$$\min_{\mathbf{a}_{0,...,N-1}} \mathbb{E}\left[\left.T(\mathbf{s}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k}, \mathbf{a}_{k})\,\right|\,\mathbf{s}_{0} = \mathbf{s}\right]$$

i.e. restrict policies to fixed $a_{0,...,N-1}$

- 4 ロ ト 4 園 ト 4 園 ト 4 園 ト 9 Q G

Infinite horizon & discounted

$$oldsymbol{\pi}_{\infty}^{\star} = rg\min_{oldsymbol{\pi}} \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{a}_k)
ight]$$

Policy π : state \rightarrow action belongs to a function space

Let's transform an MDP into an MPC and understand the approximations we make

Finite-horizon equivalent:

$$\boldsymbol{\pi}_{0,\dots,N-1}^{\star} = \operatorname*{arg\,min}_{\boldsymbol{\pi}_{0,\dots,N-1}} \mathbb{E}\left[T(\mathbf{s}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k},\mathbf{a}_{k})\right]$$

If
$$T=V^\star$$
, then $oldsymbol{\pi}_{0,\dots,N-1}^\star=oldsymbol{\pi}_\infty^\star$

Planning instead of a policy:

$$\min_{\mathbf{a}_{0,\ldots,N-1}} \mathbb{E}\left[\left.T(\mathbf{s}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k}, \mathbf{a}_{k})\right| \mathbf{s}_{0} = \mathbf{s}\right] \qquad \min_{\boldsymbol{\pi}_{0,\ldots,N-1}} T(\mathbf{s}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k}, \mathbf{a}_{k})$$

i.e. restrict policies to fixed $a_{0,...,N-1}$

Deterministic model, policy

$$\min_{\boldsymbol{\pi}_{0,\ldots,N-1}} T(\mathbf{s}_N) + \sum_{k=0}^{N-1} \gamma^k L(\mathbf{s}_k, \mathbf{a}_k)$$

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

8/29

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

i.e. adopt deterministic model

Infinite horizon & discounted

$$oldsymbol{\pi}_{\infty}^{\star} = rg\min_{oldsymbol{\pi}} \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{a}_k)
ight]$$

Policy π : state \rightarrow action belongs to a function space

Let's transform an MDP into an MPC and understand the approximations we make

Finite-horizon equivalent:

$$\boldsymbol{\pi}_{0,\dots,N-1}^{\star} = \mathop{\text{arg\,min}}_{\boldsymbol{\pi}_{0,\dots,N-1}} \mathbb{E}\left[T(\mathbf{s}_{\textit{N}}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k},\mathbf{a}_{k})\right]$$

If $T = V^*$, then $\pi_{0,\dots,N-1}^* = \pi_{\infty}^*$

Why attacking the problem in these ways?

Planning instead of a policy:

$$\min_{\mathbf{a}_{0,\ldots,N-1}} \mathbb{E}\left[\left.T(\mathbf{s}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k}, \mathbf{a}_{k})\right| \mathbf{s}_{0} = \mathbf{s}\right] \qquad \min_{\mathbf{a}_{0,\ldots,N-1}} T(\mathbf{s}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{s}_{k}, \mathbf{a}_{k})$$

i.e. restrict policies to fixed $a_{0,...,N-1}$

Deterministic model, planning

$$\min_{\mathbf{a}_0,\ldots,N-1} T(\mathbf{s}_N) + \sum_{k=0}^{N-1} \gamma^k L(\mathbf{s}_k,\mathbf{a}_k)$$

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

8/29

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

i.e. adopt deterministic model

Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad T(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t
$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k})$$

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

$$oldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

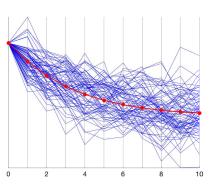
How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

Infinite horizon & discounted

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



4□ > 4□ > 4□ > 4□ > 4□ > 900

9/29

Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad T(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t
$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k})$$

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

$$oldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

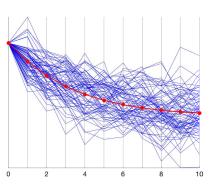
How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

Infinite horizon & discounted

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



4□ > 4□ > 4□ > 4□ > 4□ > 900

9/29

Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$

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

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

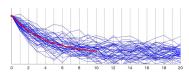
$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

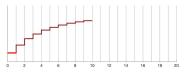
How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

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





Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$

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

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

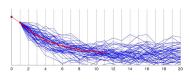
$$oldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

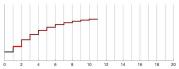
How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

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





Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$

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

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

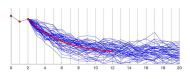
$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

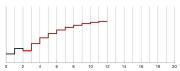
How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

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





Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$

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

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

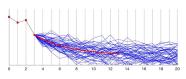
$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

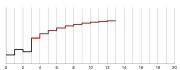
How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

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





Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t
$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k})$$

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

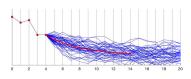
How does π^{MPC} relate to π^* ?

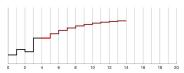
No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

Infinite horizon & discounted

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





S. Gros (NTNU)

Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t
$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k})$$

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

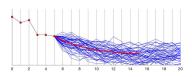
How does π^{MPC} relate to π^{\star} ?

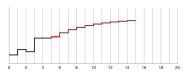
No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

Infinite horizon & discounted

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





S. Gros (NTNU)

Intro to RL-MF

Fall 2025

Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$

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

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

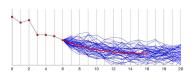
$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

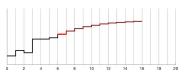
How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

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





Deterministic MPC:

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

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

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

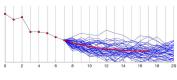
How does π^{MPC} relate to π^* ?

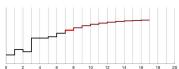
No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

Infinite horizon & discounted

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





S. Gros (NTNU)

Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t
$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k})$$

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

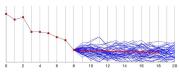
How does π^{MPC} relate to π^* ?

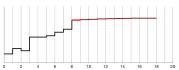
No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

Infinite horizon & discounted

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





S. Gros (NTNU)

Intro to RL-MF

Fall 2025

Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t
$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k})$$

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

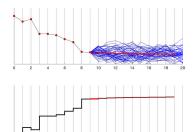
$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

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



Deterministic MPC:

$$\min_{\mathbf{u}_{0,...,N-1}} \quad \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t
$$\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k})$$

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

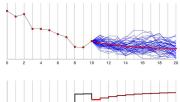
$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

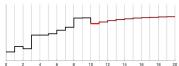
How does π^{MPC} relate to π^{\star} ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

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





Deterministic MPC:

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

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

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

$$oldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

How does π^{MPC} relate to π^{\star} ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

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



Deterministic MPC:

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

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

$$\mathbf{x}_{0} = \mathbf{s}$$

Defines policy:

$$\boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$$

How does π^{MPC} relate to π^* ?

No reason to match:

- Planning rather than policing
- Model approximates stochasticity, often deterministic

Infinite horizon & discounted

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



Can we clarify the relationship?

Historically MPC focuses on **constraints satisfaction & stability**. Cost is for reference tracking, not representative of the system performance.

S. Gros (NTNU) Intro to RL-MPC Fall 2025 10 / 29

Historically MPC focuses on **constraints satisfaction & stability**. Cost is for reference tracking, not representative of the system performance.

- "Tracking MPC"
- Classic stability theory
- Uncertainty via
 - Robust MPC
 - Stochastic MPC
- "MPC is for constraints satisfaction" (undisclosed speaker)

10 / 29

Historically MPC focuses on **constraints satisfaction & stability**. Cost is for reference tracking, not representative of the system performance.

More recently, focus shifted to closed-loop performance, e.g. energy, time, money. Cost is generic, representative of the system performance.

- "Tracking MPC"
- Classic stability theory
- Uncertainty via
 - Robust MPC
 - Stochastic MPC
- "MPC is for constraints satisfaction" (undisclosed speaker)

Historically MPC focuses on **constraints satisfaction & stability**. Cost is for reference tracking, not representative of the system performance.

- "Tracking MPC"
- Classic stability theory
- Uncertainty via
 - Robust MPC
 - Stochastic MPC
- "MPC is for constraints satisfaction" (undisclosed speaker)

More recently, focus shifted to closed-loop performance, e.g. energy, time, money. Cost is generic, representative of the system performance.

- "Economic MPC"
- Dissipativity theory
- Uncertainty via
 - Robust MPC
 - Stochastic MPC
- MPC can optimize the system performance...

Historically MPC focuses on **constraints satisfaction & stability**. Cost is for reference tracking, not representative of the system performance.

- "Tracking MPC"
- Classic stability theory
- Uncertainty via
 - Robust MPC
 - Stochastic MPC
- "MPC is for constraints satisfaction" (undisclosed speaker)

MPC for closed-loop performance

- is not a very old topic
- "historical robustness" of MPC does not hold in the presence of stochasticity

More recently, focus shifted to closed-loop performance, e.g. energy, time, money. Cost is generic, representative of the system performance.

- "Economic MPC"
- Dissipativity theory
- Uncertainty via
 - Robust MPC
 - Stochastic MPC
- MPC can optimize the system performance...

10 / 29

Historically MPC focuses on **constraints satisfaction & stability**. Cost is for reference tracking, not representative of the system performance.

- "Tracking MPC"
- Classic stability theory
- Uncertainty via
 - Robust MPC
 - Stochastic MPC
- "MPC is for constraints satisfaction" (undisclosed speaker)

MPC for closed-loop performance

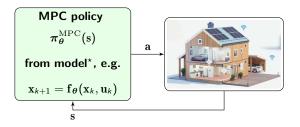
- is not a very old topic
- "historical robustness" of MPC does not hold in the presence of stochasticity

More recently, focus shifted to closed-loop performance, e.g. energy, time, money. Cost is generic, representative of the system performance.

- "Economic MPC"
- Dissipativity theory
- Uncertainty via
 - Robust MPC
 - Stochastic MPC
- MPC can optimize the system performance...

Soft claim: MPC can be used as a practical toolbox to model the solution of MDPs. This view is the "best way" for understanding what we are doing with economic MPC

Outline Let's rebuild some background - MPC and MDP MPC as a solution to MDPs



S. Gros (NTNU) Intro to RL-MPC Fall 2025 12 / 29



12 / 29

S. Gros (NTNU) Intro to RL-MPC Fall 2025



"Machine-Learning" in-the-loop f_{θ} from

- Physics-based: first principles + SYSID
- Neural Network: DNN, LSTM, TFT, . . .
- Statistical: GP, GPC, Probabilistic Al . . .

12/29

S. Gros (NTNU)



"Machine-Learning" in-the-loop f_{θ} from

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

^{*}can replace "model" by any prediction strategies: input-output predictors, multi-step predictors, etc...



"Machine-Learning" in-the-loop f_{θ} from

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

*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

12 / 29

S. Gros (NTNU) Intro to RL-MPC Fall 2025



"Machine-Learning" in-the-loop f_θ from

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

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...

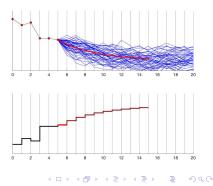
We focus on "breaking" this paradigm
Learning / RL plays a key role

MDP optimal policy

$$oldsymbol{\pi}^{\star} = \mathop{\mathsf{arg\,min}}_{oldsymbol{\pi}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{u}_k)
ight]$$

is entirely described by $Q^{\star}\left(\mathbf{s},\mathbf{a}\right)$

$$\begin{aligned} & \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star} \text{ from} \\ & \underset{\mathbf{x},\mathbf{u}}{\min} \quad \boldsymbol{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} \boldsymbol{L}(\mathbf{x}_{k},\mathbf{u}_{k}) \\ & \text{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k},\mathbf{u}_{k}), \quad \mathbf{x}_{0} = \mathbf{s} \\ & \quad \mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq \mathbf{0} \end{aligned}$$



MDP optimal policy

$$oldsymbol{\pi}^{\star} = rg\min_{oldsymbol{\pi}} \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{u}_k)
ight]$$

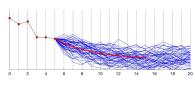
is entirely described by $Q^{\star}\left(\mathbf{s},\mathbf{a}\right)$

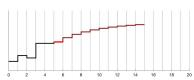
MPC as a model of the MDP

$$\begin{split} V^{\mathrm{MPC}}\left(\mathbf{s}\right) &:= \min_{\mathbf{x},\mathbf{u}} \ \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k},\mathbf{u}_{k}) \\ &\mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k},\mathbf{u}_{k}), \ \mathbf{x}_{0} = \mathbf{s} \\ &\mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq 0 \end{split}$$

$$\begin{split} \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(s\right) &= \mathbf{u}_{0}^{\star} \text{ from} \\ \min_{\mathbf{x},\mathbf{u}} \quad \boldsymbol{T}(\mathbf{x}_{\textit{N}}) + \sum_{k=0}^{\textit{N}-1} \boldsymbol{\gamma}^{k} \boldsymbol{L}(\mathbf{x}_{\textit{k}},\mathbf{u}_{\textit{k}}) \\ \mathrm{s.t.} \quad \mathbf{x}_{\textit{k}+1} &= \mathbf{f}(\mathbf{x}_{\textit{k}},\mathbf{u}_{\textit{k}}), \quad \mathbf{x}_{0} = \mathbf{s} \end{split}$$

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





MDP optimal policy

$$\boldsymbol{\pi}^{\star} = \mathop{\mathsf{arg\,min}}_{\boldsymbol{\pi}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{u}_k) \right]$$

is entirely described by $Q^{\star}(\mathbf{s}, \mathbf{a})$

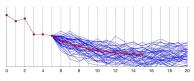
$$\begin{split} \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(s\right) &= \mathbf{u}_{0}^{\star} \text{ from} \\ \min_{\mathbf{x},\mathbf{u}} \quad \mathcal{T}(\mathbf{x}_{\textit{N}}) + \sum_{k=0}^{\textit{N}-1} \gamma^{k} \textit{L}(\mathbf{x}_{\textit{k}},\mathbf{u}_{\textit{k}}) \\ \mathrm{s.t} \quad \mathbf{x}_{\textit{k}+1} &= \mathbf{f}(\mathbf{x}_{\textit{k}},\mathbf{u}_{\textit{k}}), \quad \mathbf{x}_{0} = \mathbf{s} \\ \mathbf{h}\left(\mathbf{x}_{\textit{k}},\mathbf{u}_{\textit{k}}\right) &\leq 0 \end{split}$$

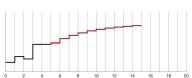
MPC as a model of the MDP

$$V^{\text{MPC}}(\mathbf{s}) := \min_{\mathbf{x}, \mathbf{u}} T(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t $\mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k}), \mathbf{x}_{0} = \mathbf{s}$

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

- Solving MDP \leftarrow building model of Q^*
- E.g. Q-learning does that from data
- MPC can model a value function...
- Can it model an action-value function?





MDP optimal policy

$$oldsymbol{\pi}^{\star} = \mathop{\mathsf{arg\,min}}_{oldsymbol{\pi}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{u}_k)
ight]$$

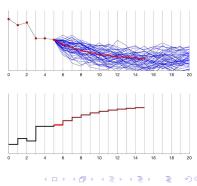
is entirely described by $Q^{\star}\left(\mathbf{s},\mathbf{a}\right)$

MPC as a model of the MDP

$$\begin{split} V^{\mathrm{MPC}}\left(\mathbf{s}\right) &:= \min_{\mathbf{x},\mathbf{u}} \ T(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k},\mathbf{u}_{k}) \\ &\mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k},\mathbf{u}_{k}), \ \mathbf{x}_{0} = \mathbf{s} \\ &\mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq 0 \end{split}$$

$$\begin{split} Q^{\mathrm{MPC}}\left(\mathbf{s},\mathbf{a}\right) &:= \min_{\mathbf{x},\mathbf{u}} \ \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} \mathcal{L}(\mathbf{x}_{k},\mathbf{u}_{k}) \\ &\mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k},\mathbf{u}_{k}) \\ &\quad \mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq 0 \\ &\quad \mathbf{x}_{0} = \mathbf{s}, \quad \mathbf{u}_{0} = \mathbf{a} \end{split}$$

$\begin{aligned} & \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(s\right) = \mathbf{u}_{0}^{\star} \text{ from} \\ & \underset{\mathbf{x},\mathbf{u}}{\min} \quad \boldsymbol{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} \boldsymbol{L}(\mathbf{x}_{k},\mathbf{u}_{k}) \\ & \text{s.t} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k},\mathbf{u}_{k}), \quad \mathbf{x}_{0} = \mathbf{s} \\ & \quad \mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq \mathbf{0} \end{aligned}$



MDP optimal policy

$$oldsymbol{\pi}^{\star} = rg\min_{oldsymbol{\pi}} \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{u}_k)
ight]$$

is entirely described by $Q^{\star}(\mathbf{s}, \mathbf{a})$

MPC as a model of the MDP

$$\begin{split} V^{\mathrm{MPC}}\left(\mathbf{s}\right) &:= \min_{\mathbf{x}, \mathbf{u}} \ \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} \mathcal{L}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ &\mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k}), \ \mathbf{x}_{0} = \mathbf{s} \\ &\mathbf{h}\left(\mathbf{x}_{k}, \mathbf{u}_{k}\right) \leq 0 \end{split}$$

$$\begin{split} Q^{\mathrm{MPC}}\left(\mathbf{s},\mathbf{a}\right) &:= \min_{\mathbf{x},\mathbf{u}} \ \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} \mathcal{L}(\mathbf{x}_{k},\mathbf{u}_{k}) \\ &\mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k},\mathbf{u}_{k}) \\ & \quad \quad \mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq 0 \\ & \quad \quad \mathbf{x}_{0} = \mathbf{s}, \quad \mathbf{u}_{0} = \mathbf{a} \end{split}$$

MPC policy $\pi^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star}$ from $\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} L(\mathbf{x}_{k}, \mathbf{u}_{k})$ $\mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k}), \quad \mathbf{x}_{0} = \mathbf{s}$ $\mathbf{h}\left(\mathbf{x}_{k}, \mathbf{u}_{k}\right) \leq \mathbf{0}$

MPC is consistent (for correct T):

$$egin{aligned} V^{\mathrm{MPC}}\left(\mathbf{s}
ight) &= \min_{\mathbf{a}} \; Q^{\mathrm{MPC}}\left(\mathbf{s},\mathbf{a}
ight) \ \pi^{\mathrm{MPC}}\left(\mathbf{s}
ight) &= rg \min_{\mathbf{a}} \; Q^{\mathrm{MPC}}\left(\mathbf{s},\mathbf{a}
ight) \end{aligned}$$

MDP optimal policy

$$m{\pi}^{\star} = \mathop{\mathsf{arg\,min}}_{m{\pi}} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k \mathit{L}(\mathbf{x}_k, \mathbf{u}_k) \right]$$

is entirely described by $Q^{\star}\left(\mathbf{s},\mathbf{a}\right)$

MPC as a model of the MDP

$$\begin{split} V^{\mathrm{MPC}}\left(\mathbf{s}\right) &:= \min_{\mathbf{x}, \mathbf{u}} \ \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} \mathcal{L}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ &\mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k}), \ \mathbf{x}_{0} = \mathbf{s} \\ &\mathbf{h}\left(\mathbf{x}_{k}, \mathbf{u}_{k}\right) \leq 0 \end{split}$$

$$\begin{split} Q^{\mathrm{MPC}}\left(\mathbf{s},\mathbf{a}\right) := \min_{\mathbf{x},\mathbf{u}} \ \mathcal{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \gamma^{k} \mathit{L}(\mathbf{x}_{k},\mathbf{u}_{k}) \\ \mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k},\mathbf{u}_{k}) \\ \quad \mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq 0 \\ \quad \mathbf{x}_{0} = \mathbf{s}, \quad \mathbf{u}_{0} = \mathbf{a} \end{split}$$

MPC policy $\pi^{\mathrm{MPC}}\left(s\right)=u_{0}^{\star}$ from

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

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

MPC is consistent (for correct T):

$$egin{aligned} V^{\mathrm{MPC}}\left(\mathbf{s}
ight) &= \min_{\mathbf{a}} \; Q^{\mathrm{MPC}}\left(\mathbf{s},\mathbf{a}
ight) \ \pi^{\mathrm{MPC}}\left(\mathbf{s}
ight) &= rg \min_{\mathbf{a}} \; Q^{\mathrm{MPC}}\left(\mathbf{s},\mathbf{a}
ight) \end{aligned}$$

MPC is a complete model of MDP if:

$$Q^{ ext{MPC}}\left(\mathbf{s},\mathbf{a}
ight)=Q^{\star}\left(\mathbf{s},\mathbf{a}
ight)$$

for all \mathbf{s}, \mathbf{a} . Then **optimality** holds:

$$oldsymbol{\pi}^{ ext{MPC}}\left(ext{s}
ight)=oldsymbol{\pi}^{\star}\left(ext{s}
ight)$$

S. Gros (NTNU)

Shift 1: focus on performance instead of fitting

- from: f_{θ} is a model for the system dynamics
- to: MPC is a model of the MDP solution

S. Gros (NTNU) Intro to RL-MPC Fall 2025 14 / 29

Shift 1: focus on performance instead of fitting

- from: f_{θ} is a model for the system dynamics
- to: MPC is a model of the MDP solution

Classic view...

$$\begin{split} \textbf{MPC: at current state s solve} \\ \min_{\mathbf{x},\mathbf{u}} \quad & \mathcal{T}\left(\mathbf{x}_{\textit{N}}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{\textit{k}},\mathbf{u}_{\textit{k}}\right) \\ \text{s.t.} \quad & \mathbf{x}_{k+1} = \mathbf{f}_{\theta}\left(\mathbf{x}_{\textit{k}},\mathbf{u}_{\textit{k}}\right) \\ & \quad & \mathbf{h}\left(\mathbf{x}_{\textit{k}},\mathbf{u}_{\textit{k}}\right) \leq 0 \\ & \quad & \mathbf{x}_{0} = \mathbf{s} \end{split}$$
 gives policy $\boldsymbol{\pi}_{\theta}^{\text{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star}$

Find θ such that prediction "fits" the data

S. Gros (NTNU) Intro to RL-MP

Shift 1: focus on performance instead of fitting

- from: f_{θ} is a model for the system dynamics
- to: MPC is a model of the MDP solution

Classic view...

$$\begin{aligned} & \text{MPC: at current state } \mathbf{s} \text{ solve} \\ & \min_{\mathbf{x}, \mathbf{u}} \quad \mathcal{T}\left(\mathbf{x}_N\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_k, \mathbf{u}_k\right) \\ & \text{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}_{\theta}\left(\mathbf{x}_k, \mathbf{u}_k\right) \\ & \quad \mathbf{h}\left(\mathbf{x}_k, \mathbf{u}_k\right) \leq 0 \\ & \quad \mathbf{x}_0 = \mathbf{s} \end{aligned}$$
 gives policy $\boldsymbol{\pi}_{\theta}^{\text{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star}$

Find θ such that prediction "fits" the data

Shift to...

Find θ that "fits MPC to optimality" according to the data $(Q^{\mathrm{MPC}} o Q^{\star} ext{ or at least } \pi^{\mathrm{MPC}} o \pi^{\star})$

Shift 1: focus on performance instead of fitting

- from: f_{θ} is a model for the system dynamics
- to: MPC is a model of the MDP solution

Classic view...

$$\begin{aligned} \textbf{MPC:} & \text{ at current state } \mathbf{s} \text{ solve} \\ & \min_{\mathbf{x},\mathbf{u}} \quad \mathcal{T}\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \text{ s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \quad \mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq 0 \\ & \quad \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$

Find θ such that prediction "fits" the data

gives policy $\boldsymbol{\pi}_{\theta}^{\mathrm{MPC}}(\mathbf{s}) = \mathbf{u}_{0}^{\star}$

Shift to...

Find θ that "fits MPC to optimality" according to the data

$$\left(\mathit{Q}^{\mathrm{MPC}}
ightarrow \mathit{Q}^{\star}$$
 or at least $oldsymbol{\pi}^{\mathrm{MPC}}
ightarrow oldsymbol{\pi}^{\star}
ight)$

- ullet Best model for closed-loop performance
- $\bullet \neq \mathsf{Best} \mathsf{ model to fit the data!}$

S. Gros (NTNU) Intro to RL-MPC

Shift 1: focus on performance instead of fitting

- from: f_{θ} is a model for the system dynamics
- to: MPC is a model of the MDP solution

Classic view...

MPC: at current state s solve

$$\min_{x \in \mathcal{X}_{N}} T(x_{N}) + \sum_{x \in \mathcal{X}_{N}} I(x_{N}, y_{N})$$

$$\min_{\mathbf{x},\mathbf{u}} \quad T\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{\infty} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$$

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
$$oldsymbol{\pi}_{ heta}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

Find θ such that prediction "fits" the data

Shift to...

Find θ that "fits MPC to optimality" according to the data

$$(\mathit{Q}^{\mathrm{MPC}}
ightarrow \mathit{Q}^{\star}$$
 or at least $oldsymbol{\pi}^{\mathrm{MPC}}
ightarrow oldsymbol{\pi}^{\star})$

- $\bullet \ \to \mathsf{Best} \ \mathsf{model} \ \mathsf{for} \ \mathsf{closed}\text{-loop} \ \mathsf{performance}$
- $\bullet \neq \mathsf{Best} \mathsf{ model to fit the data!}$

But
$$oldsymbol{\pi}^{\mathrm{MPC}} = oldsymbol{\pi}^{\star}$$
 places "high demands" on $\mathbf{f}_{ heta}$

Can we do better?

14 / 29

S. Gros (NTNU) Intro to RL-MPC Fall 2025

Shift 1: focus on performance instead of fitting

- from: f_{θ} is a model for the system dynamics
- to: MPC is a model of the MDP solution

Classic view...

MPC: at current state s solve
$$\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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
$$oldsymbol{\pi}_{ heta}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$

Find θ such that prediction "fits" the data

Shift to...

Find θ that "fits MPC to optimality" according to the data

(
$$Q^{\mathrm{MPC}}
ightarrow Q^\star$$
 or at least $m{\pi}^{\mathrm{MPC}}
ightarrow m{\pi}^\star$)

- $\bullet \ \to \mathsf{Best} \ \mathsf{model} \ \mathsf{for} \ \mathsf{closed}\text{-loop} \ \mathsf{performance}$
- $\bullet \neq \mathsf{Best} \mathsf{ model to fit the data!}$

But
$$oldsymbol{\pi}^{\mathrm{MPC}} = oldsymbol{\pi}^{\star}$$
 places "high demands" on $\mathbf{f}_{ heta}$

Can we do better? Yes!

14 / 29

4□ > 4□ > 4 亘 > 4 亘 > □ ● 9 Q @

Shift 1: focus on performance instead of fitting

- from: f_{θ} is a model for the system dynamics
- to: MPC is a model of the MDP solution

Classic view...

$$\begin{aligned} & \text{MPC: at current state } \mathbf{s} \text{ solve} \\ & \min_{\mathbf{x},\mathbf{u}} \quad \mathcal{T}\left(\mathbf{x}_{\textit{N}}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \text{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \quad \mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq 0 \\ & \quad \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$

Find θ such that prediction "fits" the data

gives policy $\pi_{\theta}^{\mathrm{MPC}}(\mathbf{s}) = \mathbf{u}_{0}^{\star}$

Shift to...

Find θ that "fits MPC to optimality" according to the data

$$(\mathit{Q}^{\mathrm{MPC}}
ightarrow \mathit{Q}^{\star}$$
 or at least $oldsymbol{\pi}^{\mathrm{MPC}}
ightarrow oldsymbol{\pi}^{\star})$

Shift 2: "holistic" parametrization

$$\min_{\mathbf{x},\mathbf{u}} \quad T_{\theta}\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{N-1} L_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$$

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

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

i.e. cost and constraints are part of the model

Shift 1: focus on performance instead of fitting

- from: f_{θ} is a model for the system dynamics
- to: MPC is a model of the MDP solution

Classic view...

$$\begin{split} \textbf{MPC:} & \text{ at current state } \textbf{s} \text{ solve} \\ & \min_{\textbf{x},\textbf{u}} \quad \mathcal{T}\left(\textbf{x}_{\textit{N}}\right) + \sum_{k=0}^{N-1} L\left(\textbf{x}_{\textit{k}},\textbf{u}_{\textit{k}}\right) \\ & \text{s.t.} \quad \textbf{x}_{\textit{k}+1} = \textbf{f}_{\theta}\left(\textbf{x}_{\textit{k}},\textbf{u}_{\textit{k}}\right) \\ & \quad \textbf{h}\left(\textbf{x}_{\textit{k}},\textbf{u}_{\textit{k}}\right) \leq 0 \\ & \quad \textbf{x}_{0} = \textbf{s} \\ & \text{gives policy } \boldsymbol{\pi}_{\theta}^{\text{MPC}}\left(\textbf{s}\right) = \textbf{u}_{0}^{\star} \end{split}$$

Find θ such that prediction "fits" the data

Shift to...

Find θ that "fits MPC to optimality" according to the data $(Q^{\mathrm{MPC}} \to Q^{\star} \text{ or at least } \pi^{\mathrm{MPC}} \to \pi^{\star})$

Full MPC parametrization:

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

$$& \text{gives policy } \boldsymbol{\pi}_{\theta}^{\text{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star}$$

S. Gros (NTNU) Intro to RL-MPC Fall 2025 15 / 29

Full MPC parametrization:

$$egin{aligned} & \min_{\mathbf{x},\mathbf{u}} & T_{ heta}\left(\mathbf{x}_{N}
ight) + \sum_{k=0}^{N-1} L_{ heta}\left(\mathbf{x}_{k},\mathbf{u}_{k}
ight) \\ & ext{s.t.} & \mathbf{x}_{k+1} = \mathbf{f}_{ heta}\left(\mathbf{x}_{k},\mathbf{u}_{k}
ight) \\ & & \mathbf{h}_{ heta}\left(\mathbf{x}_{k},\mathbf{u}_{k}
ight) \leq 0 \\ & & \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$
 gives policy $oldsymbol{\pi}_{ heta}^{ ext{MPC}}\left(\mathbf{s}
ight) = \mathbf{u}_{0}^{\star}$

Theorem: under some technical conditions and for a "rich" parametrization of the MPC, there is a θ such that

$$egin{aligned} Q_{ heta}^{ ext{MPC}}\left(\mathbf{s},\mathbf{a}
ight) &= Q^{\star}\left(\mathbf{s},\mathbf{a}
ight) \ oldsymbol{\pi}_{ heta}^{ ext{MPC}}\left(\mathbf{s}
ight) &= oldsymbol{\pi}^{\star}\left(\mathbf{s}
ight) \end{aligned}$$

even if the model \mathbf{f}_{θ} cannot describe the real system accurately

S. Gros (NTNU) Intro to RL-MPC Fall 2025 15 / 29

Full MPC parametrization:

$$\min_{\mathbf{x},\mathbf{u}} \quad T_{\theta}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} L_{\theta}(\mathbf{x}_{k}, \mathbf{u}_{k})$$
s.t.
$$\mathbf{x}_{k+1} = \mathbf{f}_{\theta}(\mathbf{x}_{k}, \mathbf{u}_{k})$$

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

$$\mathbf{x}_{0} = \mathbf{s}$$

gives policy $\pi_{\theta}^{\mathrm{MPC}}(\mathbf{s}) = \mathbf{u}_{0}^{\star}$

Theorem: under some technical conditions and for a "rich" parametrization of the MPC, there is a θ such that

$$egin{aligned} Q_{ heta}^{ ext{MPC}}\left(ext{s}, ext{a}
ight) &= Q^{\star}\left(ext{s}, ext{a}
ight) \ oldsymbol{\pi}_{ heta}^{ ext{MPC}}\left(ext{s}
ight) &= oldsymbol{\pi}^{\star}\left(ext{s}
ight) \end{aligned}$$

even if the model $\mathbf{f}_{ heta}$ cannot describe the real system accurately

Remarks

 Compensate MPC model deficiencies in the cost + constraints

S. Gros (NTNU) Intro to RL-MPC

Full MPC parametrization:

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

gives policy $\boldsymbol{\pi}_{\theta}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star}$

Theorem: under some technical conditions and for a "rich" parametrization of the MPC, there is a θ such that

$$egin{aligned} Q_{ heta}^{ ext{MPC}}\left(ext{s}, ext{a}
ight) &= Q^{\star}\left(ext{s}, ext{a}
ight) \ \pi_{ heta}^{ ext{MPC}}\left(ext{s}
ight) &= \pi^{\star}\left(ext{s}
ight) \end{aligned}$$

even if the model \mathbf{f}_{θ} cannot describe the real system accurately

Remarks

- Compensate MPC model deficiencies in the cost + constraints
- Generic: works for robust MPC, stochastic MPC, economic MPC, Multi-Step PC, etc.

S. Gros (NTNU) Intro to RL-MPC Fall 2025 15 / 29

Full MPC parametrization:

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

gives policy $\boldsymbol{\pi}_{\theta}^{\mathrm{MPC}}\left(\mathbf{s}\right)=\mathbf{u}_{0}^{\star}$

Theorem: under some technical conditions and for a "rich" parametrization of the MPC, there is a θ such that

$$egin{aligned} Q_{ heta}^{ ext{MPC}}\left(ext{s}, ext{a}
ight) &= Q^{\star}\left(ext{s}, ext{a}
ight) \ \pi_{ heta}^{ ext{MPC}}\left(ext{s}
ight) &= \pi^{\star}\left(ext{s}
ight) \end{aligned}$$

even if the model \mathbf{f}_{θ} cannot describe the real system accurately

Remarks

- Compensate MPC model deficiencies in the cost + constraints
- Generic: works for robust MPC, stochastic MPC, economic MPC, Multi-Step PC, etc.
- Sanity check: technical conditions are mild but forbid MPC model to be "very absurd"

Full MPC parametrization:

$$egin{aligned} \min_{\mathbf{x},\mathbf{u}} \quad & T_{ heta}\left(\mathbf{x}_{N}
ight) + \sum_{k=0}^{N-1} L_{ heta}\left(\mathbf{x}_{k},\mathbf{u}_{k}
ight) \\ \mathrm{s.t.} \quad & \mathbf{x}_{k+1} = \mathbf{f}_{ heta}\left(\mathbf{x}_{k},\mathbf{u}_{k}
ight) \\ & \mathbf{h}_{ heta}\left(\mathbf{x}_{k},\mathbf{u}_{k}
ight) \leq 0 \\ & \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$
 gives policy $\pi_{ heta}^{\mathrm{MPC}}\left(\mathbf{s}
ight) = \mathbf{u}_{0}^{\star}$

MPC is a model of the optimal policy

not a policy approximation using open-loop (model-based) predictions

Theorem: under some technical conditions and for a "rich" parametrization of the MPC, there is a θ such that

$$egin{aligned} Q_{ heta}^{ ext{MPC}}\left(ext{s}, ext{a}
ight) &= Q^{\star}\left(ext{s}, ext{a}
ight) \ oldsymbol{\pi}_{ heta}^{ ext{MPC}}\left(ext{s}
ight) &= oldsymbol{\pi}^{\star}\left(ext{s}
ight) \end{aligned}$$

even if the model f_{θ} cannot describe the real system accurately

Remarks

- Compensate MPC model deficiencies in the cost + constraints
- Generic: works for robust MPC, stochastic MPC, economic MPC, Multi-Step PC, etc.
- Sanity check: technical conditions are mild but forbid MPC model to be "very absurd"

□ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 □ >
 <l

How to use this? Reinforcement Learning

$$\begin{split} & \text{Policy } \boldsymbol{\pi}_{\theta}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star} \text{ from} \\ & \underset{\mathbf{x},\mathbf{u}}{\min} \quad \boldsymbol{T}_{\theta}\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{N-1} \boldsymbol{L}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \text{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \quad \quad \mathbf{h}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq \mathbf{0}, \quad \mathbf{x}_{0} = \mathbf{s} \end{split}$$

- ullet min $_{ heta} J\left(oldsymbol{\pi}_{ heta}^{ ext{MPC}}
 ight)$ using data
- $m{ heta}$ $heta o J\left(m{\pi}_{ heta}^{ ext{MPC}}
 ight)$ very implicit
- J(.) is the real-system!

◄□▶◀圖▶◀불▶◀불▶ 불 쒸٩€

16 / 29

S. Gros (NTNU) Intro to RL-MPC Fall 2025

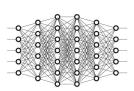
How to use this? Reinforcement Learning

$$\begin{aligned} & \text{Policy } \boldsymbol{\pi}_{\theta}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star} \text{ from} \\ & \underset{\mathbf{x},\mathbf{u}}{\min} \quad \boldsymbol{\mathcal{T}}_{\theta}\left(\mathbf{x}_{\textit{N}}\right) + \sum_{k=0}^{N-1} \boldsymbol{\mathcal{L}}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \text{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \quad \mathbf{h}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq \mathbf{0}, \quad \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$

- ullet min $_{ heta} J\left(oldsymbol{\pi}_{ heta}^{ ext{MPC}}
 ight)$ using data
- ullet $heta o J\left(oldsymbol{\pi}_{ heta}^{\mathrm{MPC}}
 ight)$ very implicit
- J(.) is the real-system!

Reinforcement Learning

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





How to use this? Reinforcement Learning

Policy
$$oldsymbol{\pi}_{ heta}^{\mathrm{MPC}}\left(\mathbf{s}
ight)=\mathbf{u}_{0}^{\star}$$
 from

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

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

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

- ullet min $_{ heta} J\left(oldsymbol{\pi}_{ heta}^{ ext{MPC}}
 ight)$ using data
- $oldsymbol{ heta}$ $heta o J\left(oldsymbol{\pi}_{ heta}^{\mathrm{MPC}}
 ight)$ very implicit
- J(.) is the real-system!

Reinforcement Learning

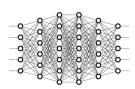
Tools to approximate π^* from data This is not (necessarily) about DNNs

For MPC: tools to find best θ , e.g.

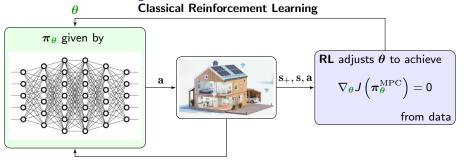
Policy Gradient: estimations of

$$abla_{ heta} J\left(m{\pi}_{ heta}^{\mathrm{MPC}}
ight), \quad ext{possibly} \quad
abla_{ heta}^2 J\left(m{\pi}_{ heta}^{\mathrm{MPC}}
ight)$$

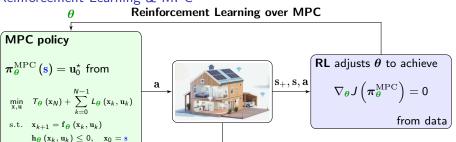
• Q-learning: direct "shaping" of MPC







S. Gros (NTNU)

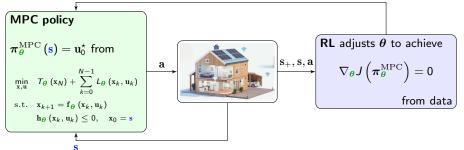


S. Gros (NTNU)

Intro to RL-MPC

Fall 2025

 θ Reinforcement Learning over MPC

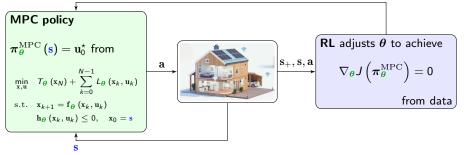


Why is it useful?

- ✓ Tunes your optimization model for real-world performance
- √ MPC: 25 years of results on formal guarantees & methods
- √ Easy to inject knowledge
- √ Learning does not start from scratch
- √ Planning provides explainability
- ✓ Software now available



 θ Reinforcement Learning over MPC



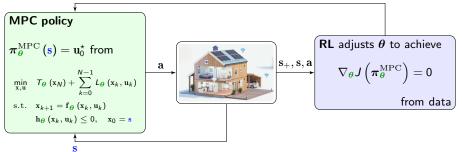
Why is it useful?

- ✓ Tunes your optimization model for real-world performance
- √ MPC: 25 years of results on formal guarantees & methods
- √ Easy to inject knowledge
- √ Learning does not start from scratch
- √ Planning provides explainability
- √ Software now available

Why can it be difficult?

- Optimization is computationally more expensive than a DNN
- ✗ Deployment on GPUs is lagging
- Non-convexity can get in the way...





Why is it useful?

- ✓ Tunes your optimization model for real-world performance
- √ MPC: 25 years of results on formal guarantees & methods
- √ Easy to inject knowledge
- √ Learning does not start from scratch
- √ Planning provides explainability
- √ Software now available

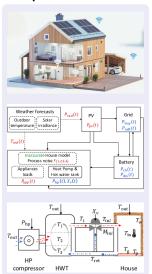
Why can it be difficult?

- Optimization is computationally more expensive than a DNN
- ✗ Deployment on GPUs is lagging
- X Non-convexity can get in the way...

RL over MPC can tune your MPC for performance beyond what classical methods can do!

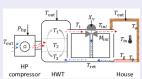
Example - Home Energy Management (simulated)

Setup



Setup

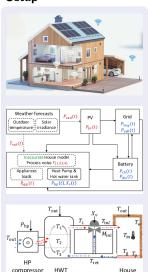




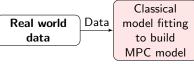
Learning process: aligned followed with closed-loop

Real world data

Setup



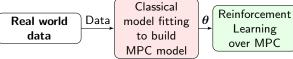
Learning process: aligned followed with closed-loop



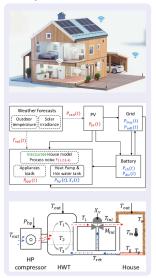
Setup



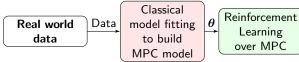
Learning process: aligned followed with closed-loop



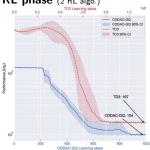
Setup



Learning process: aligned followed with closed-loop



RL phase (2 RL algo.)



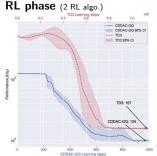
Setup



Real world data

Data Classical model fitting to build to build over MPC

MPC model



The RL step improves the MPC performance significantly from only model fitting

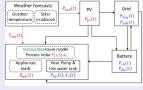
HWT

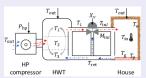
House

compressor

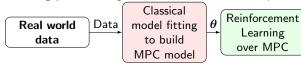
Setup



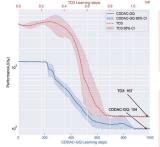




Learning process: aligned followed with closed-loop



RL phase (2 RL algo.)



The RL step improves the MPC performance significantly from only model fitting

RL over MPC can tune your MPC for performance beyond what classical methods can do!

Outline Let's rebuild some background - MPC and MDP When is RL (most) beneficial for MPC?

Model fitting?

Proposed paradigm

Policy
$$\pi_{\theta}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star}$$
 from
$$\min_{\mathbf{x},\mathbf{u}} \quad T_{\theta}\left(\mathbf{x}_{\mathit{N}}\right) + \sum_{k=0}^{\mathit{N}-1} L_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$$
 s.t. $\mathbf{x}_{k+1} = \mathbf{f}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$ $\mathbf{h}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq 0, \quad \mathbf{x}_{0} = \mathbf{s}$

S. Gros (NTNU) Intro to RL-MPC Fall 2025 20 / 29

Model fitting?

A step back: model adjustment?

Policy
$$\boldsymbol{\pi}_{\theta}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star}$$
 from
$$\begin{aligned} & \min_{\mathbf{x},\mathbf{u}} \quad \boldsymbol{T}\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \\ & \quad \mathbf{h}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right) \leq \mathbf{0}, \quad \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$

Adjust θ according to (e.g.)

- $1. \quad \mathop{\text{min}}_{\theta} \mathbb{E}\left[\left\|\mathbf{f}_{\theta}\left(s,a\right)-s_{+}\right\|^{2}\right] \text{ vs. }$
- 2. $\min_{\theta} J\left(\boldsymbol{\pi}_{\theta}^{\text{MPC}}\right)$

... is that different?

S. Gros (NTNU)

Intro to RL-MPC

Model fitting?

A step back: model adjustment?

Policy
$$oldsymbol{\pi}^{\mathrm{MPC}}_{ heta}\left(\mathbf{s}\right) = \mathbf{u}^{\star}_{0} \; \mathsf{from}$$

$$\min_{\mathbf{x},\mathbf{u}} \;\; \mathcal{T}\left(\mathbf{x}_{\mathit{N}}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$$

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

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

Adjust θ according to (e.g.)

- $1. \quad \mathop{\text{min}}_{\theta} \mathbb{E}\left[\left\| \mathbf{f}_{\theta}\left(\mathbf{s}, \mathbf{a}\right) \mathbf{s}_{+} \right\|^{2} \right] \text{ vs.}$
- 2. $\min_{\theta} J\left(\boldsymbol{\pi}_{\theta}^{\mathrm{MPC}}\right)$

... is that different?

Empirical answers:

- In general it is different...
- Good to do 1. first, then 2.
- Is step 2 necessary?
 - Improves closed-loop performance
 - But gain is *very* case dependent (also with full parametrization)

What causes that?

S. Gros (NTNU)

Intro to RL-MF

$$\begin{aligned} & \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_0^{\star} \text{ from} \\ & \min_{\mathbf{u}, \mathbf{x}} \quad \mathcal{T}(\mathbf{x}_{\mathit{N}}) + \sum_{k=0}^{N-1} \mathcal{L}(\mathbf{x}_{\mathit{k}}, \mathbf{u}_{\mathit{k}}) \\ & \text{s.t} \quad \mathbf{x}_{\mathit{k}+1} = \mathbf{f}(\mathbf{x}_{\mathit{k}}, \mathbf{u}_{\mathit{k}}) \\ & \quad \mathbf{h}(\mathbf{x}_{\mathit{k}}, \mathbf{u}_{\mathit{k}}) \leq 0, \quad \mathbf{x}_0 = \mathbf{s} \end{aligned}$$

S. Gros (NTNU) Intro to RL-MPC Fall 2025 21 / 29

$$\begin{split} \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(s\right) &= \mathbf{u}_0^{\star} \text{ from} \\ \min_{\mathbf{u}, \mathbf{x}} \quad \boldsymbol{T}(\mathbf{x}_{\textit{N}}) + \sum_{k=0}^{\textit{N}-1} \boldsymbol{L}(\mathbf{x}_{\textit{k}}, \mathbf{u}_{\textit{k}}) \\ \mathrm{s.t} \quad \mathbf{x}_{\textit{k}+1} &= \mathbf{f}(\mathbf{x}_{\textit{k}}, \mathbf{u}_{\textit{k}}) \\ \mathbf{h}(\mathbf{x}_{\textit{k}}, \mathbf{u}_{\textit{k}}) &\leq 0, \quad \mathbf{x}_0 = \mathbf{s} \end{split}$$

Sufficient condition of optimality of MPC relates model ${\bf f}$ to optimal value function and conditional distribution ${\bf s}, {\bf a} \to {\bf s}_+$ of the MDP (+ condition on T)

S. Gros (NTNU) Intro to RL-MPC Fall 2025 21 / 29

$$\begin{split} \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(s\right) &= \mathbf{u}_0^{\star} \text{ from} \\ \min_{\mathbf{u}, \mathbf{x}} \quad \boldsymbol{T}(\mathbf{x}_{\textit{N}}) + \sum_{k=0}^{\textit{N}-1} \boldsymbol{L}(\mathbf{x}_{\textit{k}}, \mathbf{u}_{\textit{k}}) \\ \mathrm{s.t} \quad \mathbf{x}_{k+1} &= \mathbf{f}(\mathbf{x}_{\textit{k}}, \mathbf{u}_{\textit{k}}) \\ \mathbf{h}(\mathbf{x}_{\textit{k}}, \mathbf{u}_{\textit{k}}) \leq 0, \quad \mathbf{x}_0 &= \mathbf{s} \end{split}$$

Sufficient condition of optimality of MPC relates model ${\bf f}$ to optimal value function and conditional distribution ${\bf s}, {\bf a} \to {\bf s}_+$ of the MDP (+ condition on T)

In general

ridge regression

$$f\left(s,a\right)=\mathbb{E}\left[\,s_{+}\,|\,s,a\,\right]$$

max likelihood

$$\mathbf{f}\left(\mathbf{s}, \mathbf{a}\right) = \mathsf{max}\, \rho\left[\,\mathbf{s}_{+} \,|\, \mathbf{s}, \mathbf{a}\,
ight]$$

do not satisfy the optimality conditions

$$\begin{aligned} & \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(s\right) = \mathbf{u}_{0}^{\star} \text{ from} \\ & \underset{\mathbf{u}, \mathbf{x}}{\min} \quad \boldsymbol{T}(\mathbf{x}_{\textit{N}}) + \sum_{k=0}^{\textit{N}-1} \boldsymbol{L}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ & \text{s.t} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ & \quad \mathbf{h}(\mathbf{x}_{k}, \mathbf{u}_{k}) < 0, \quad \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$

Sufficient condition of optimality of MPC relates model ${\bf f}$ to optimal value function and conditional distribution ${\bf s}, {\bf a} \to {\bf s}_+$ of the MDP (+ condition on T)

In general

ridge regression

$$f\left(s,a\right)=\mathbb{E}\left[\,s_{+}\,|\,s,a\,\right]$$

max likelihood

$$\mathbf{f}\left(\mathbf{s}, \mathbf{a}\right) = \mathsf{max}\, \rho\left[\,\mathbf{s}_{+} \,|\, \mathbf{s}, \mathbf{a}\,
ight]$$

do not satisfy the optimality conditions

Model from "classic SYSID" does not necessarily yield the best MPC policy

S. Gros (NTNU)

Intro to RL-MPC

Fall 2025

$$\begin{aligned} \textbf{MPC policy } & \boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star} \text{ from} \\ & \min_{\mathbf{u}, \mathbf{x}} \quad \boldsymbol{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} \boldsymbol{L}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ & \text{s.t} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ & \quad \mathbf{h}(\mathbf{x}_{k}, \mathbf{u}_{k}) \leq 0, \quad \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$

Sufficient condition of optimality of MPC relates model ${\bf f}$ to optimal value function and conditional distribution ${\bf s}, {\bf a} \to {\bf s}_+$ of the MDP (+ condition on T)

In general

ridge regression

$$f\left(s,a\right)=\mathbb{E}\left[\,s_{+}\,|\,s,a\,\right]$$

max likelihood

$$f(s, a) = \max \rho[s_+ | s, a]$$

do not satisfy the optimality conditions

Model from "classic SYSID" does not necessarily yield the best MPC policy

Notable exceptions: model gives $\mathbb{E}\left[\,\mathbf{s}_{+}\,|\,\mathbf{s},\mathbf{a}\,\right]$ and

- LQR with process noise
- By extension, locally optimal policies for
 - dissipative MDP and
 - $ightharpoonup V^\starpprox$ quadratic in positive invariant set

Smooth tracking MPC, fixed reference away from constraints, "reasonable" stochasticity

$$\begin{aligned} & \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right) = \mathbf{u}_{0}^{\star} \text{ from} \\ & \min_{\mathbf{u}, \mathbf{x}} \quad \boldsymbol{T}(\mathbf{x}_{N}) + \sum_{k=0}^{N-1} L(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ & \text{s.t} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ & \quad \mathbf{h}(\mathbf{x}_{k}, \mathbf{u}_{k}) < 0, \quad \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$

Sufficient condition of optimality of MPC relates model ${\bf f}$ to optimal value function and conditional distribution ${\bf s}, {\bf a} \to {\bf s}_+$ of the MDP (+ condition on T)

When is RL (most) useful for MPC? "practical" view

 Not much if smooth tracking problem, spend most of the time close to fixed steady state, rare transients away from the constraints → small performance gain

21/29

$$\begin{aligned} & \text{MPC policy } \boldsymbol{\pi}^{\mathrm{MPC}}\left(s\right) = \mathbf{u}_{0}^{\star} \text{ from} \\ & \underset{\mathbf{u}, \mathbf{x}}{\min} \quad \boldsymbol{T}(\mathbf{x}_{\textit{N}}) + \sum_{k=0}^{\textit{N}-1} \boldsymbol{L}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ & \text{s.t} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_{k}, \mathbf{u}_{k}) \\ & \quad \mathbf{h}(\mathbf{x}_{\textit{k}}, \mathbf{u}_{\textit{k}}) < 0, \quad \mathbf{x}_{0} = \mathbf{s} \end{aligned}$$

Sufficient condition of optimality of MPC relates model ${\bf f}$ to optimal value function and conditional distribution ${\bf s}, {\bf a} \to {\bf s}_+$ of the MDP (+ condition on T)

When is RL (most) useful for MPC? "practical" view

- Not much if smooth tracking problem, spend most of the time close to fixed steady state, rare transients away from the constraints → small performance gain
- Model is not rich enough to predict the (relevant) expected state transition
- Economic problem with "low dissipativity" (trajectories spread over the state space)
- Value function curvature changes a lot, problem is non-smooth
- Optimal steady-state near / at constraints
- Varying exogeneous inputs (e.g. varying prices in energy-related problems, changing reference)
- Task-based problems: racing, minimum time, termination conditions, etc.

4□ > 4□ > 4≡ > 4≡ > 900

21/29

S. Gros (NTNU) Intro to RL-MPC Fall 2025

Outline Let's rebuild some background - MPC and MDP A Deeper Look at the Theory

The theory is about equivalences between MDPs

S. Gros (NTNU)

The theory is about equivalences between MDPs

World MDP

World MDP definition

- States s and actions a
- Cost *L*(s, a)
- Transition $s_+ \sim \rho \left[\cdot | s, a \right]$

The theory is about equivalences between MDPs

World MDP

World MDP definition

- States s and actions a
- Cost *L*(s, a)
- Transition $\mathbf{s}_+ \sim \rho \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Optimal value functions

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

$$Q^{\star}\left(\mathbf{s},\mathbf{a}
ight)=\mathbf{\emph{L}}\left(\mathbf{s},\mathbf{a}
ight)+\gamma\mathbb{E}_{\mathbf{s}_{+}\simoldsymbol{
ho}}\left[V^{\star}\left(\mathbf{s}_{+}
ight)|\mathbf{s},\mathbf{a}
ight]$$

The theory is about equivalences between MDPs

World MDP

World MDP definition

- States s and actions a
- Cost *L*(s, a)
- Transition $\mathbf{s}_+ \sim \rho \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Optimal value functions

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

$$\mathbf{\textit{Q}}^{\star}\left(\mathbf{s},\mathbf{a}\right)=\textcolor{red}{\textit{L}}\left(\mathbf{s},\mathbf{a}\right)+\gamma\mathbb{E}_{\mathbf{s}_{+}\sim\textcolor{red}{\rho}}\left[\textit{V}^{\star}\left(\mathbf{s}_{+}\right)|\mathbf{s},\mathbf{a}\right]$$

Optimal policy

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

The theory is about equivalences between MDPs

World MDP

World MDP definition

- States s and actions a
- Cost *L*(s, a)
- Transition $\mathbf{s}_+ \sim \rho \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Optimal value functions

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

$$\mathbf{Q}^{\star}\left(\mathbf{s},\mathbf{a}
ight) = \mathbf{L}\left(\mathbf{s},\mathbf{a}
ight) + \gamma \mathbb{E}_{\mathbf{s}_{+} \sim \mathbf{\rho}}\left[\mathbf{V}^{\star}\left(\mathbf{s}_{+}
ight) | \mathbf{s},\mathbf{a}
ight]$$

Optimal policy

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

Model MDP

Model MDP definition

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\pmb{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Optimal value functions

$$\hat{V}^{\star}\left(\mathbf{s}\right) = \min_{\mathbf{a}} \; \hat{Q}^{\star}\left(\mathbf{s}, \mathbf{a}\right)$$

$$\hat{Q}^{\star}\left(\mathbf{s},\mathbf{a}\right) = \hat{L}\left(\mathbf{s},\mathbf{a}\right) + \gamma \mathbb{E}_{\mathbf{s}_{+} \sim \hat{\rho}}\left[\left|V^{\star}\left(\hat{\mathbf{s}}_{+}\right)\right|\left|\mathbf{s},\mathbf{a}\right]\right]$$

Optimal policy

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

The theory is about equivalences between MDPs

World MDP

World MDP definition

- States s and actions a
- Cost *L*(s, a)
- Transition $\mathbf{s}_+ \sim \rho \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Optimal value functions

$$V^{\star}\left(\mathbf{s}\right) = \min_{\mathbf{a}} \ Q^{\star}\left(\mathbf{s}, \mathbf{a}\right)$$

$$Q^{\star}\left(\mathbf{s},\mathbf{a}\right) = \frac{L}{L}\left(\mathbf{s},\mathbf{a}\right) + \gamma \mathbb{E}_{\mathbf{s}_{+} \sim \rho}\left[V^{\star}\left(\mathbf{s}_{+}\right) | \mathbf{s}, \mathbf{a}\right]$$

Optimal policy

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

Model MDP

Model MDP definition

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\rho} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Optimal value functions

$$\hat{V}^{\star}\left(\mathbf{s}\right) = \min_{\mathbf{a}} \; \hat{Q}^{\star}\left(\mathbf{s}, \mathbf{a}\right)$$

$$\hat{Q}^{\star}\left(\mathbf{s},\mathbf{a}\right) = \hat{L}\left(\mathbf{s},\mathbf{a}\right) + \gamma \mathbb{E}_{\mathbf{s}_{+} \sim \hat{\rho}}\left[\left|V^{\star}\left(\hat{\mathbf{s}}_{+}\right)\right|\left|\mathbf{s},\mathbf{a}\right]\right]$$

Optimal policy

$$\hat{\pi}^{\star}(\mathbf{s}) = \underset{\mathbf{a}}{\mathsf{arg\,min}} \ \hat{Q}^{\star}(\mathbf{s}, \mathbf{a})$$

Theory says that - under some technical conditions - there is a \hat{L} such that $\hat{Q}^\star = Q^\star$

The theory is about equivalences between MDPs

World MDP

World MDP definition

- States s and actions a
- Cost *L*(s, a)
- Transition $s_+ \sim \rho \left[\cdot | s, a \right]$

Optimal value functions

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

$$\mathbf{\textit{Q}}^{\star}\left(\mathbf{s},\mathbf{a}\right)=\textcolor{red}{\textit{L}}\left(\mathbf{s},\mathbf{a}\right)+\gamma\mathbb{E}_{\mathbf{s}_{+}\sim\textcolor{red}{\rho}}\left[\textit{V}^{\star}\left(\mathbf{s}_{+}\right)|\mathbf{s},\mathbf{a}\right]$$

Optimal policy

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

Model MDP

Model MDP definition

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\rho} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Optimal value functions

$$\hat{V}^{\star}\left(\mathbf{s}
ight)=\min_{\mathbf{a}}\;\hat{Q}^{\star}\left(\mathbf{s},\mathbf{a}
ight)$$

$$\hat{Q}^{\star}\left(\mathbf{s},\mathbf{a}\right) = \hat{L}\left(\mathbf{s},\mathbf{a}\right) + \gamma \mathbb{E}_{\mathbf{s}_{+} \sim \hat{\rho}}\left[\left|V^{\star}\left(\hat{\mathbf{s}}_{+}\right)\right|\left|\mathbf{s},\mathbf{a}\right]\right]$$

Optimal policy

$$\hat{\pi}^{\star}\left(\mathbf{s}\right) = \underset{\mathbf{a}}{\operatorname{arg\,min}} \ \hat{Q}^{\star}\left(\mathbf{s}, \mathbf{a}\right)$$

Theory says that - under some technical conditions - there is a $\hat{m L}$ such that $\hat{Q}^\star = Q^\star$

 $\underline{\mathsf{Proof}}\!\!: \text{ telescopic sum, some non-trivial assumptions to prevent } \infty - \infty \text{ cancellations}$

S. Gros (NTNU) Intro to RL-MPC Fall 2025 23 / 29

World MDP

- States s and actions a
- \bullet Cost L(s, a)
- Transition $\mathbf{s}_+ \sim \rho \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Model MDP

- States s and actions a
 - Cost $\hat{L}(s, a)$
 - Transition $\mathbf{s}_+ \sim \hat{\boldsymbol{\rho}} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Theory says that

ullet Under some technical conditions - there is a $\hat{\mathcal{L}}$ such that $\hat{Q}^\star = Q^\star$

World MDP

- States s and actions a
- Cost *L*(s, a)
- Transition $\mathbf{s}_+ \sim \frac{\rho}{\rho} [\cdot | \mathbf{s}, \mathbf{a}]$

Model MDP

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\boldsymbol{\rho}} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- ullet For $\hat{L}={\color{red}L}$ there is a set of models $\hat{
 ho}$ such that $\hat{Q}^{\star}={\color{red}Q}^{\star}$

World MDP

- States s and actions a
- Cost *L*(s, a)
- Transition $\mathbf{s}_+ \sim \frac{\rho}{\rho} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Model MDP

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\boldsymbol{\rho}} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^{\star} = Q^{\star}$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)

World MDP

- States s and actions a
- Cost *L*(s, a)
- Transition $\mathbf{s}_+ \sim \frac{\rho}{\rho} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Model MDP

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\boldsymbol{\rho}} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Theory says that

- ullet Under some technical conditions there is a $\hat{\mathcal{L}}$ such that $\hat{\mathcal{Q}}^\star = \mathcal{Q}^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^{\star} = Q^{\star}$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)
- ullet World MDP and Model MDP do not need to use the same discount γ

S. Gros (NTNU)

World MDP

- States s and actions a
- Cost **L**(s, a)
- ullet Transition $\mathbf{s}_+ \sim oldsymbol{
 ho} \left[\, \cdot \, | \mathbf{s}, \mathbf{a}
 ight]$

Model MDP

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\pmb{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^{\star} = Q^{\star}$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)
- ullet World MDP and Model MDP do not need to use the same discount γ

MPC is a "weird" undiscounted Model MDP

- Deterministic model is a trivial $\hat{\rho}$, then optimal policy for Model MDP \equiv planning
- Finite horizon: ok if correct terminal cost (strong assumption, but learning can tune it)
- Scenario tree: rough $\hat{\rho}$ and policy.
- Robust MPC: carries support $\hat{\rho}$, worst-case cost is a bit off...

World MDP

- States s and actions a
- Cost *L*(s, a)
- ullet Transition $\mathbf{s}_+ \sim oldsymbol{
 ho} \left[\, \cdot \, | \mathbf{s}, \mathbf{a}
 ight]$

Model MDP

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\boldsymbol{\rho}} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Theory says that

- ullet Under some technical conditions there is a $\hat{\mathcal{L}}$ such that $\hat{\mathcal{Q}}^\star = \mathcal{Q}^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^{\star} = Q^{\star}$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)
- ullet World MDP and Model MDP do not need to use the same discount γ

Remarks

• Non-unique optimal model $\hat{\rho}$ leaves room for aligning it

World MDP

- States s and actions a
- Cost *L*(s, a)
- Transition $\mathbf{s}_+ \sim \frac{\rho}{\rho} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Model MDP

- States s and actions a
- Cost $\hat{L}(s, a)$
- Transition $\mathbf{s}_+ \sim \hat{\boldsymbol{\rho}} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^{\star} = Q^{\star}$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)
- ullet World MDP and Model MDP do not need to use the same discount γ

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then $f(s,a)\subset \text{support of }\rho$

World MDP

- States s and actions a
- Cost **L**(s, a)
- Transition $\mathbf{s}_+ \sim \frac{\rho}{\rho} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Model MDP

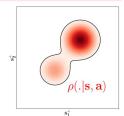
- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^* = Q^*$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)
- ullet World MDP and Model MDP do not need to use the same discount γ

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then $f(s,a)\subset \text{support of }\rho$



24 / 29

S. Gros (NTNU) Intro to RL-MPC Fall 2025

World MDP

- States s and actions a
- Cost **L**(s, a)
- Transition $\mathbf{s}_+ \sim \frac{\rho}{\rho} \left[\cdot | \mathbf{s}, \mathbf{a} \right]$

Model MDP

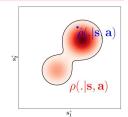
- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^* = Q^*$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)
- ullet World MDP and Model MDP do not need to use the same discount γ

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then $f(s,a) \subset \text{support of } \rho$



24 / 29

S. Gros (NTNU) Intro to RL-MPC Fall 2025

World MDP

- States s and actions a
- Cost **L**(s, a)
- ullet Transition $s_+ \sim {\color{red}
 ho} \left[\, \cdot \, | s, a
 ight]$

Model MDP

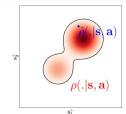
- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^{\star} = Q^{\star}$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)
- lacktriangle World MDP and Model MDP do not need to use the same discount γ

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then $\mathbf{f}(\mathbf{s},\mathbf{a})\subset \text{support of }\rho$
- More simply said: we can build a "plausible" optimal deterministic model f (predictions have prob. > 0).



4 T > 4 A > 4 B > 4 B > B > 9 Q

More on MDP Equivalences

World MDP

- States s and actions a
- Cost **L**(s, a)
- Transition $\mathbf{s}_+ \sim \rho \left[\, \cdot \, | \mathbf{s}, \mathbf{a} \right]$

Model MDP

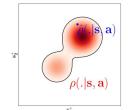
- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- For $\hat{L} = L$ there is a set of models $\hat{\rho}$ such that $\hat{Q}^{\star} = Q^{\star}$
- Conditions for model $\hat{\rho}$ "optimality" \neq min of classical loss functions (except. LQR)
- ullet World MDP and Model MDP do not need to use the same discount γ

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then $\mathbf{f}(\mathbf{s},\mathbf{a})\subset \text{support of }\rho$
- More simply said: we can build a "plausible" optimal deterministic model f (predictions have prob. > 0).
- Cost modification promotes "simplicity".



World MDP

- States s and actions a
- Cost L (s, a)
- \bullet Transition $s_+ \sim \frac{\rho}{\rho} \left[\, \cdot \, | s, a \right]$

Model MDP

- States s and actions a
- Cost $\hat{L}(\mathbf{s}, \mathbf{a})$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\, \cdot \, | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a $\hat{m L}$ such that $\hat{Q}^\star = Q^\star$

Technical assumptions? They are very technical and not intuitive at all. But not very demanding, and there is a "self-fulfilling" effect in learning, i.e. RL "learns" to satisfy them.

World MDP

- States s and actions a
- Cost L (s, a)
- $\bullet \ \, \mathsf{Transition} \,\, \mathbf{s}_{+} \sim \textcolor{red}{\rho} \, [\, \cdot \, | \mathbf{s}, \mathbf{a}]$

Model MDP

- States s and actions a
- Cost $\hat{L}(\mathbf{s}, \mathbf{a})$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\, \cdot \, | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a $\hat{m L}$ such that $\hat{Q}^\star = Q^\star$

But we are making a strong assumption (in plain sight)

S. Gros (NTNU)

Intro to RL-MPC

World MDP

- States s and actions a
- Cost L (s, a)
- ullet Transition $s_+ \sim {\color{red}
 ho} \left[\, \cdot \, | s, a
 ight]$

Model MDP

- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{
 ho} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a \hat{L} such that $\hat{Q}^\star = Q^\star$

But we are making a strong assumption (in plain sight)

- \bullet Theory assumes that World MDP and Model MDP have the same state s
- I.e. we need to know the state of the World MDP to build a correct Model MDP
- That is often unrealistic...

World MDP

- States s and actions a
- Cost *L*(s, a)
- $\bullet \ \ \mathsf{Transition} \ \mathbf{s}_{+} \sim {\color{red} \rho} \left[\, \cdot \, | \mathbf{s}, \mathbf{a} \right]$

Model MDP

- States s and actions a
- Cost $\hat{L}(\mathbf{s}, \mathbf{a})$
- ullet Transition $\mathbf{s}_+ \sim \hat{
 ho} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a \hat{L} such that $\hat{Q}^\star = Q^\star$

But we are making a strong assumption (in plain sight)

- \bullet Theory assumes that World MDP and Model MDP have the same state s
- I.e. we need to know the state of the World MDP to build a correct Model MDP
- That is often unrealistic...

The end of a nice theory? No, but it is a word of caution...

World MDP

- States s and actions a
- Cost L (s, a)
- Transition $s_+ \sim \rho \left[\cdot | s, a \right]$

Model MDP

- States s and actions a
- Cost $\hat{L}(\mathbf{s}, \mathbf{a})$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\, \cdot \, | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a \hat{L} such that $\hat{Q}^\star = Q^\star$

Eventually the Markov state s "boils down to" the history of past observations

How to embed that in a Model MDP?

World MDP

- States s and actions a
- Cost L (s, a)
- ullet Transition $s_+ \sim {\color{red}
 ho} \left[\, \cdot \, | s, a
 ight]$

Model MDP

- States s and actions a
- Cost $\hat{L}(\mathbf{s}, \mathbf{a})$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\, \cdot \, | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a \hat{L} such that $\hat{Q}^\star = Q^\star$

Eventually the Markov state ${\bf s}$ "boils down to" the history of past observations How to embed that in a Model MDP?

1 "Input-output models": ARX, multi-step predictors (linear or not), etc.

World MDP

- States s and actions a
- Cost L (s, a)
- $\bullet \ \, \mathsf{Transition} \,\, \mathbf{s}_{+} \sim \textcolor{red}{\rho} \, [\, \cdot \, | \mathbf{s}, \mathbf{a}]$

Model MDP

- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{
 ho} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a \hat{L} such that $\hat{Q}^\star = Q^\star$

Eventually the Markov state ${\bf s}$ "boils down to" the history of past observations How to embed that in a Model MDP?

- 1 "Input-output models": ARX, multi-step predictors (linear or not), etc.
- 2 Latent states (a.k.a. embeddings or compressed representations) giving an Al-driven lower-dimensional version of that history

World MDP

- States s and actions a
- Cost L (s, a)
- ullet Transition $s_+ \sim {\color{red}
 ho} \left[\, \cdot \, | s, a
 ight]$

Model MDP

- States s and actions a
- Cost $\hat{L}(\mathbf{s}, \mathbf{a})$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\, \cdot \, | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a \hat{L} such that $\hat{Q}^\star = Q^\star$

Eventually the Markov state s "boils down to" the history of past observations How to embed that in a Model MDP?

- 1 "Input-output models": ARX, multi-step predictors (linear or not), etc.
- 2 Latent states (a.k.a. embeddings or compressed representations) giving an Al-driven lower-dimensional version of that history
- State observers giving a model-based meaningful low-dimensional version of that history

World MDP

- States s and actions a
- Cost *L*(s, a)
- ullet Transition $\mathbf{s}_+ \sim oldsymbol{
 ho} \left[\, \cdot \, | \mathbf{s}, \mathbf{a}
 ight]$

Model MDP

- States s and actions a
- Cost $\hat{L}(\mathbf{s}, \mathbf{a})$
- ullet Transition $\mathbf{s}_+ \sim \hat{
 ho} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that - under some technical conditions - there is a $\hat{m L}$ such that $\hat{Q}^\star = Q^\star$

Eventually the Markov state s "boils down to" the history of past observations

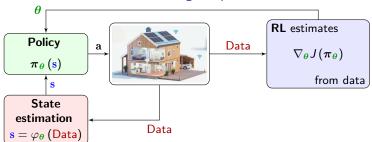
How to embed that in a Model MDP?

- 1 "Input-output models": ARX, multi-step predictors (linear or not), etc.
- 2 Latent states (a.k.a. embeddings or compressed representations) giving an Al-driven lower-dimensional version of that history
- State observers giving a model-based meaningful low-dimensional version of that history

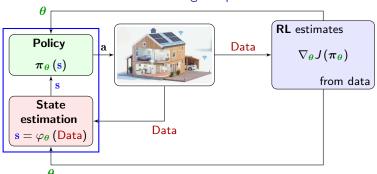
Options 2 and 3 are part of the decision-making process, back-propagate through policy + state "constructor".

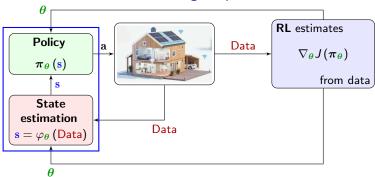
S. Gros (NTNU) Intro to RL-MPC Fall 2025 25 / 29





S. Gros (NTNU) Intro to RL-MPC Fall 2025 26 / 29

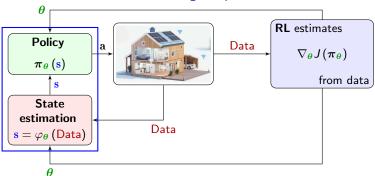




Remarks

ullet Policy is $oldsymbol{\pi}_{ heta}^{ extsf{D}}\left(extsf{Data}
ight)=oldsymbol{\pi}_{ heta}\left(arphi_{ heta}\left(extsf{Data}
ight)
ight)$

S. Gros (NTNU) Intro to RL-MPC Fall 2025 26 / 29

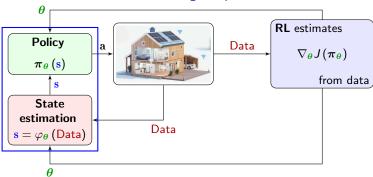


Remarks

- Policy is $\pi_{\theta}^{\mathsf{D}}(\mathsf{Data}) = \pi_{\theta}(\varphi_{\theta}(\mathsf{Data}))$
- Gradient of the policy now is:

$$\nabla_{\theta} \boldsymbol{\pi}^{\mathsf{D}}_{\theta} \left(\mathsf{Data} \right) = \nabla_{\theta} \boldsymbol{\pi}_{\theta} \left(\mathbf{s} \right) + \nabla_{\theta} \varphi_{\theta} \left(\mathsf{Data} \right) \nabla_{\mathbf{s}} \boldsymbol{\pi}_{\theta} \left(\mathbf{s} \right)$$

◆ロト ◆個ト ◆差ト ◆差ト を めるぐ



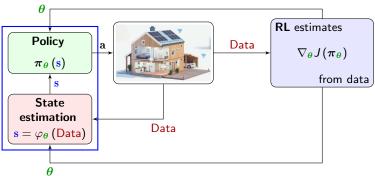
Remarks

- Policy is $\pi_{\theta}^{\mathsf{D}}(\mathsf{Data}) = \pi_{\theta}(\varphi_{\theta}(\mathsf{Data}))$
- Gradient of the policy now is:

$$\nabla_{\theta} \boldsymbol{\pi}_{\theta}^{\mathsf{D}} \left(\mathsf{Data} \right) = \nabla_{\theta} \boldsymbol{\pi}_{\theta} \left(\mathbf{s} \right) + \nabla_{\theta} \varphi_{\theta} \left(\mathsf{Data} \right) \nabla_{\mathbf{s}} \boldsymbol{\pi}_{\theta} \left(\mathbf{s} \right)$$

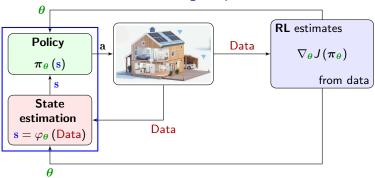
• Critic is looking at Data $\xrightarrow{\theta}$ a, i.e. $Q^{\pi_{\theta}^{D}}$ (Data, a)

◆ロト ◆御 ト ◆ 恵 ト ◆ 恵 ・ 釣 ९ ○



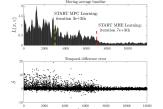
If φ_{θ} (Data) is

- MHE, then its gradient is computed using NLP sensitivities (same as MPC)
- DNN, then efficient sensitivity computations exist
- Bayesian inference? Belief state? Open research topic...



If φ_{θ} (Data) is

- MHE, then its gradient is computed using NLP sensitivities (same as MPC)
- DNN, then efficient sensitivity computations exist
- Bayesian inference? Belief state? Open research topic...



26 / 29

Systems with

- ∼Linear dynamics
- Input-output data
- Significant stochasticity
- Modelling is difficult

27 / 29

S. Gros (NTNU) Intro to RL-MPC Fall 2025

Systems with

- ◆ Linear dynamics
- Input-output data
- Significant stochasticity
- Modelling is difficult

Multi-step linear predictors

$$\hat{\mathbf{y}} = \Phi \left[egin{array}{c} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{array}
ight]$$

- Recent history of input-output sequence u, y
- ullet Planned input sequence ${\bf u}$
- Predicted output sequence \hat{y}

S. Gros (NTNU) Intro to RL-MPC Fall 2025 27 / 29

SPC for u, y given

$$\min_{\mathbf{u},\,\hat{\mathbf{y}}} \quad \sum_{k=0}^{N} L(\hat{\mathbf{y}}_{k},\mathbf{u}_{k})$$

s.t.
$$\hat{\mathbf{y}} = \Phi \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

 $\mathbf{h}(\hat{\mathbf{y}}_k, \mathbf{u}_k) < 0$

yields policy
$$\pi(\mathbf{u}, \mathbf{y}) = \mathbf{u}_0^*$$

Multi-step linear predictors

$$\hat{\mathbf{y}} = \Phi \left[egin{array}{c} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{array}
ight]$$

- \bullet Recent history of input-output sequence \mathbf{u},\mathbf{y}
- ullet Planned input sequence ${\bf u}$
- Predicted output sequence \hat{y}

SPC for u, y given

$$\min_{\mathbf{u}, \hat{\mathbf{y}}} \quad \sum_{k=0}^{N} L(\hat{\mathbf{y}}_{k}, \mathbf{u}_{k})$$
s.t.
$$\hat{\mathbf{y}} = \Phi \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

$$\mathbf{h}\left(\hat{\mathbf{y}}_{k},\mathbf{u}_{k}\right)<0$$

yields policy
$$\pi\left(\mathbf{u},\mathbf{y}\right)=\mathbf{u}_{0}^{\star}$$

Multi-step linear predictors

$$\hat{\mathbf{y}} = \Phi \left[egin{array}{c} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{array}
ight]$$

- ullet Recent history of input-output sequence \mathbf{u},\mathbf{y}
- ullet Planned input sequence ${\bf u}$
- Predicted output sequence \hat{y}
- Measured output sequences y

matrix Φ can be **built from past data** \mathcal{D} , e.g.

$$\min_{\Phi} \sum_{i \in \mathcal{D}} \frac{1}{2} \left\| \mathbf{y}_{i} - \Phi \begin{bmatrix} \mathbf{u}_{i} \\ \mathbf{y}_{i} \\ \mathbf{u}_{i} \end{bmatrix} \right\|^{2} + R(\Phi)$$

s.t. Φ causal

or alternative loss functions (e.g. quantile)

◆ロト ◆母 ト ◆ 恵 ト ◆ 恵 ・ 釣 へ ○

27 / 29

SPC for u, y given

$$\min_{\mathbf{u},\,\hat{\mathbf{y}}} \quad \sum_{k=0}^{N} L(\hat{\mathbf{y}}_{k},\mathbf{u}_{k})$$

s.t.
$$\hat{\mathbf{y}} = \Phi \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

 $\mathbf{h} (\hat{\mathbf{y}}_k, \mathbf{u}_k) < 0$

yields policy
$$\pi\left(\mathbf{u},\mathbf{y}\right) = \mathbf{u}_0^{\star}$$

If R promotes a dynamic system structure, regression well-posed even with limited data (and Φ high-dimensional)

But best model fit *⇒* best policy

Multi-step linear predictors

$$\hat{\mathbf{y}} = \Phi \left[egin{array}{c} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{array}
ight]$$

- ullet Recent history of input-output sequence \mathbf{u},\mathbf{y}
- ullet Planned input sequence ${\bf u}$
- Predicted output sequence \hat{y}
- Measured output sequences y

matrix Φ can be **built from past data** \mathcal{D} , e.g.

$$\min_{\Phi} \sum_{i \in \mathcal{D}} \frac{1}{2} \left\| \mathbf{y}_{i} - \Phi \begin{bmatrix} \mathbf{u}_{i} \\ \mathbf{y}_{i} \\ \mathbf{u}_{i} \end{bmatrix} \right\|^{2} + R(\Phi)$$

s.t. Φ causal

or alternative loss functions (e.g. quantile)

S. Gros (NTNU) Intro to RL-MPC Fall 2025 27 / 29

SPC for \mathbf{u} , \mathbf{y} given

$$\min_{\mathbf{u}, \, \hat{\mathbf{y}}} \quad \sum_{k=0}^{N} L(\hat{\mathbf{y}}_{k}, \mathbf{u}_{k})$$
s.t.
$$\hat{\mathbf{y}} = \Phi \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

yields policy
$$\pi\left(\mathbf{u},\mathbf{y}\right)=\mathbf{u}_0^{\star}$$

 $h(\hat{\mathbf{y}}_k, \mathbf{u}_k) < 0$

Can we close the gap via RL? Yes!

- RL-MPC theory applies with some twists
- State becomes u, y (window of input-output)

But best model fit $\not\Rightarrow$ best policy

27 / 29

ro to RL-MPC Fall 2025

MSPC for u, y given

$$\min_{\mathbf{u},\,\hat{\boldsymbol{y}}}\quad \Psi_{\theta}\left(\mathbf{u},\,\hat{\boldsymbol{y}},\mathbf{u},\boldsymbol{y}\right)$$

s.t.
$$\hat{\mathbf{y}} = \Phi_{\theta} \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

$$\mathbf{H}_{\theta}\left(\mathbf{u},\,\hat{\mathbf{y}},\mathbf{u},\mathbf{y}\right)\leq 0$$

yields policy
$$oldsymbol{\pi}_{oldsymbol{ heta}}\left(\mathbf{u},\mathbf{y}
ight)=\mathbf{u}_{0}^{\star}$$

Can we close the gap via RL? Yes!

- RL-MPC theory applies with some twists
- State becomes **u**, **y** (window of input-output)
- Modifications in principle not localized in time

But best model fit *⇒* best policy

Fall 2025

27 / 29

MSPC for u, y given

$$\min_{\mathbf{u},\,\hat{\boldsymbol{y}}}\quad \Psi_{\theta}\left(\mathbf{u},\,\hat{\boldsymbol{y}},\mathbf{u},\boldsymbol{y}\right)$$

s.t.
$$\hat{\mathbf{y}} = \Phi_{\theta} \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

$$\mathbf{H}_{\theta}\left(\mathbf{u},\, \boldsymbol{\mathring{\mathbf{y}}}, \boldsymbol{\overset{\mathbf{u}}{\mathbf{u}}}, \boldsymbol{\overset{\mathbf{y}}{\mathbf{y}}}\right) \leq 0$$

yields policy $oldsymbol{\pi}_{ heta}\left(\mathbf{u},\mathbf{y}
ight)=\mathbf{u}_{0}^{\star}$

Can we close the gap via RL? Yes!

- RL-MPC theory applies with some twists
- State becomes u, y (window of input-output)
- Modifications in principle not localized in time
- High-dimensional parameter space for RL

But best model fit ≠ best policy



27 / 29

tro to RL-MPC Fall 2025

MSPC for u, y given

$$\min_{\mathbf{u},\,\hat{\mathbf{y}}}\quad \Psi_{\theta}\left(\mathbf{u},\,\hat{\mathbf{y}},\mathbf{u},\mathbf{y}\right)$$

s.t.
$$\hat{\mathbf{y}} = \Phi_{\theta} \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

$$\mathbf{H}_{\theta}\left(\mathbf{u},\,\hat{\mathbf{y}},\mathbf{u},\mathbf{y}\right)\leq 0$$

yields policy $oldsymbol{\pi}_{ heta}\left(\mathbf{u},\mathbf{y}
ight)=\mathbf{u}_{0}^{\star}$

Can we close the gap via RL? Yes!

- RL-MPC theory applies with some twists
- State becomes u, y (window of input-output)
- Modifications in principle not localized in time
- High-dimensional parameter space for RL

We need an add-on in Leap-C for high-speed MPC + sensitivity on this type of model structure! For now we are down to using CVXPY.

But best model fit ≠ best policy

◆ロト ◆個ト ◆意ト ◆意ト · 恵 · のQで

27 / 29

S. Gros (NTNU) Intro to RL-MPC Fall 2025

MSPC for u, y given

$$\min_{\mathbf{u},\,\hat{\mathbf{y}}} \quad \Psi_{\theta}\left(\mathbf{u},\,\hat{\mathbf{y}},\mathbf{u},\mathbf{y}\right)$$

s.t.
$$\hat{\mathbf{y}} = \Phi_{\theta} \begin{bmatrix} \mathbf{u} \\ \mathbf{y} \\ \mathbf{u} \end{bmatrix}$$

$$\mathbf{H}_{\theta}\left(\mathbf{u},\, \boldsymbol{\mathring{\mathbf{y}}}, \mathbf{u}, \mathbf{y}\right) \leq 0$$

yields policy $oldsymbol{\pi}_{ heta}\left(\mathbf{u},\mathbf{y}
ight)=\mathbf{u}_{0}^{\star}$

Can we close the gap via RL? Yes!

- RL-MPC theory applies with some twists
- State becomes u, y (window of input-output)
- Modifications in principle not localized in time
- High-dimensional parameter space for RL

We need an add-on in Leap-C for high-speed MPC + sensitivity on this type of model structure! For now we are down to using CVXPY.

But best model fit *⇒* best policy

S. Gros (NTNU)

◆ロト ◆個ト ◆意ト ◆意ト · 恵 · のQで

27 / 29

Intro to RL-MPC Fall 2025

Orientation

What we have seen:

- MPC can be understood as a model of the optimal action-value function Q^* of real-world MDPs and/or of the optimal policy π^*
- MPC cost (and constraints) become part of that model
- Model that best fits the real-world does not (necessarily) yield the best policy
- RL is a toolbox to tune the MPC as a model of the MDP solution
- MPC state space should match the real world, strong assumption that can be alleviated

28 / 29

Orientation

What we have seen:

- MPC can be understood as a model of the optimal action-value function Q^* of real-world MDPs and/or of the optimal policy π^*
- MPC cost (and constraints) become part of that model
- Model that best fits the real-world does not (necessarily) yield the best policy
- RL is a toolbox to tune the MPC as a model of the MDP solution
- MPC state space should match the real world, strong assumption that can be alleviated

What we will do next: RL over MPC

- Safe & Stable RL over MPC (In the afternoon)
- RL over MPC with belief states a future prospect (In the afternoon)
- Beyond MPC Model-based Decisions and AI for decisions (Tomorrow)

4 D > 4 D > 4 E > 4 E > 9 Q P

28 / 29

Thanks for your attention!



ResearchGate



Google Scholar