

Differentiating MPC with applications in Reinforcement Learning

Mario Zanon

IMT School for Advanced Studies Lucca

Joint work with Sébastien Gros
ACC 2023

Iteratively solve:

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

MPC

Iteratively solve:

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

Optimal policy: $\pi_{\theta}(s) := \mathbf{u}_0^*(s, \theta)$ deterministic

MPC

Iteratively solve:

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

Optimal policy: $\pi_{\theta}(s) := \mathbf{u}_0^*(s, \theta)$ deterministic

Remarks:

Iteratively solve:

$$\mathbf{u}^*(s, \theta), \mathbf{x}^*(s, \theta) = \arg \min_{\mathbf{u}, \mathbf{x}} V_{\theta}^f(\mathbf{x}_N) + \sum_{k=0}^{N-1} \ell_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_0 = \mathbf{s},$$

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

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

$$\mathbf{h}_{\theta}^f(\mathbf{x}_N) \leq 0,$$

Optimal policy: $\pi_{\theta}(s) := \mathbf{u}_0^*(s, \theta)$ deterministic

Remarks:

- Typically
 - ▶ \mathbf{f}_{θ} “exact” or from system identification
 - ▶ $\ell_{\theta}, \mathbf{h}_{\theta}$ given
 - ▶ \mathbf{h}_{θ}^f invariant set, V_{θ}^f Lyapunov function

Iteratively solve:

$$\mathbf{u}^*(s, \theta), \mathbf{x}^*(s, \theta) = \arg \min_{\mathbf{u}, \mathbf{x}} V_{\theta}^f(\mathbf{x}_N) + \sum_{k=0}^{N-1} \ell_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_0 = \mathbf{s},$$

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

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

$$\mathbf{h}_{\theta}^f(\mathbf{x}_N) \leq 0,$$

Optimal policy: $\pi_{\theta}(s) := \mathbf{u}_0^*(s, \theta)$ deterministic

Remarks:

- Typically
 - ▶ \mathbf{f}_{θ} “exact” or from system identification
 - ▶ $\ell_{\theta}, \mathbf{h}_{\theta}$ given
 - ▶ \mathbf{h}_{θ}^f invariant set, V_{θ}^f Lyapunov function
- If \mathbf{f}_{θ} exact and $N \rightarrow \infty$: asymptotic stability + optimal closed-loop performance

Iteratively solve:

$$\mathbf{u}^*(s, \theta), \mathbf{x}^*(s, \theta) = \arg \min_{\mathbf{u}, \mathbf{x}} V_{\theta}^f(\mathbf{x}_N) + \sum_{k=0}^{N-1} \ell_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{x}_0 = \mathbf{s},$$

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

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

$$\mathbf{h}_{\theta}^f(\mathbf{x}_N) \leq 0,$$

Optimal policy: $\pi_{\theta}(s) := \mathbf{u}_0^*(s, \theta)$ deterministic

Remarks:

- Typically
 - ▶ \mathbf{f}_{θ} “exact” or from system identification
 - ▶ $\ell_{\theta}, \mathbf{h}_{\theta}$ given
 - ▶ \mathbf{h}_{θ}^f invariant set, V_{θ}^f Lyapunov function
- If \mathbf{f}_{θ} exact and $N \rightarrow \infty$: asymptotic stability + optimal closed-loop performance
- In practice: model inaccuracy \implies some form of practical stability and suboptimal

Iteratively solve:

$$\mathbf{u}^*(s, \theta), \mathbf{x}^*(s, \theta) = \arg \min_{\mathbf{u}, \mathbf{x}} V_{\theta}^f(\mathbf{x}_N) + \sum_{k=0}^{N-1} \ell_{\theta}(\mathbf{x}_k, \mathbf{u}_k)$$

s.t. $\mathbf{x}_0 = \mathbf{s}$,

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

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

$$\mathbf{h}_{\theta}^f(\mathbf{x}_N) \leq 0,$$

Optimal policy: $\pi_{\theta}(s) := \mathbf{u}_0^*(s, \theta)$ **deterministic**

Remarks:

- Typically
 - ▶ \mathbf{f}_{θ} “exact” or from system identification
 - ▶ $\ell_{\theta}, \mathbf{h}_{\theta}$ given
 - ▶ \mathbf{h}_{θ}^f invariant set, V_{θ}^f Lyapunov function
- If \mathbf{f}_{θ} exact and $N \rightarrow \infty$: asymptotic stability + optimal closed-loop performance
- In practice: model inaccuracy \implies some form of practical stability and suboptimal
- Let RL adapt θ to recover optimality

Reinforcement Learning

- A technique to solve a stochastic optimal control problem
 - ▶ Many similarities with Dynamic Programming
 - ▶ Based on sampling
 - ▶ Model based / model free

Reinforcement Learning

- A technique to solve a stochastic optimal control problem
 - ▶ Many similarities with Dynamic Programming
 - ▶ Based on sampling
 - ▶ Model based / model free
- Cross-disciplinary approach
 - ▶ Computer Science
 - ▶ Optimal Control
 - ▶ Systems / Control Theory
 - ▶ Statistics

Reinforcement Learning

- A technique to solve a stochastic optimal control problem
 - ▶ Many similarities with Dynamic Programming
 - ▶ Based on sampling
 - ▶ Model based / model free
- Cross-disciplinary approach
 - ▶ Computer Science
 - ▶ Optimal Control
 - ▶ Systems / Control Theory
 - ▶ Statistics
- Important success stories
 - ▶ Aerobatic helicopter flight
 - ▶ Win against Chess, Shogi and Go masters
 - ▶ Learn to beat the best Chess algorithm in a matter of few hours

Reinforcement Learning

- A technique to solve a stochastic optimal control problem
 - ▶ Many similarities with Dynamic Programming
 - ▶ Based on sampling
 - ▶ Model based / model free
- Cross-disciplinary approach
 - ▶ Computer Science
 - ▶ Optimal Control
 - ▶ Systems / Control Theory
 - ▶ Statistics
- Important success stories
 - ▶ Aerobatic helicopter flight
 - ▶ Win against Chess, Shogi and Go masters
 - ▶ Learn to beat the best Chess algorithm in a matter of few hours
- Limitations
 - ▶ Theoretical foundation of some algorithms not fully developed
 - ▶ **No stability and safety guarantees**

Reinforcement Learning

- A technique to solve a stochastic optimal control problem
 - ▶ Many similarities with Dynamic Programming
 - ▶ Based on sampling
 - ▶ Model based / model free
- Cross-disciplinary approach
 - ▶ Computer Science
 - ▶ Optimal Control
 - ▶ Systems / Control Theory
 - ▶ Statistics
- Important success stories
 - ▶ Aerobatic helicopter flight
 - ▶ Win against Chess, Shogi and Go masters
 - ▶ Learn to beat the best Chess algorithm in a matter of few hours
- Limitations
 - ▶ Theoretical foundation of some algorithms not fully developed
 - ▶ **No stability and safety guarantees**
- Advantages
 - ▶ **Optimality for the true system**

RL in a Nutshell

- Markov Chain defined by

$$\mathbb{P}[s_+ | s, \mathbf{a}] \quad \text{equivalent of} \quad s_+ = f(s, \mathbf{a}, \mathbf{w})$$

probability (density) of transitioning from state s to state s_+ when taking action \mathbf{a}

RL in a Nutshell

- Markov Chain defined by

$$\mathbb{P}[\mathbf{s}_+ | \mathbf{s}, \mathbf{a}] \quad \text{equivalent of} \quad \mathbf{s}_+ = \mathbf{f}(\mathbf{s}, \mathbf{a}, \mathbf{w})$$

probability (density) of transitioning from state \mathbf{s} to state \mathbf{s}_+ when taking action \mathbf{a}

- Stochastic policy (includes deterministic as special case)

$$\pi[\mathbf{a} | \mathbf{s}] \in \mathbb{R}_+ \quad \text{deterministic: } \mathbf{a} = \pi(\mathbf{s})$$

assigns the probability (density) of taking action \mathbf{a} for a given state \mathbf{s}

RL in a Nutshell

- Markov Chain defined by

$$\mathbb{P}[\mathbf{s}_+ | \mathbf{s}, \mathbf{a}] \quad \text{equivalent of} \quad \mathbf{s}_+ = \mathbf{f}(\mathbf{s}, \mathbf{a}, \mathbf{w})$$

probability (density) of transitioning from state \mathbf{s} to state \mathbf{s}_+ when taking action \mathbf{a}

- Stochastic policy (includes deterministic as special case)

$$\pi[\mathbf{a} | \mathbf{s}] \in \mathbb{R}_+ \quad \text{deterministic: } \mathbf{a} = \pi(\mathbf{s})$$

assigns the probability (density) of taking action \mathbf{a} for a given state \mathbf{s}

- Expected discounted cost (return):

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

where \mathbf{a}_k is drawn from policy π . Initial conditions \mathbf{s}_0 can be fixed or random.

- Discount factor $\gamma \in [0, 1]$

RL in a Nutshell

- Markov Chain defined by

$$\mathbb{P}[\mathbf{s}_+ | \mathbf{s}, \mathbf{a}] \quad \text{equivalent of} \quad \mathbf{s}_+ = \mathbf{f}(\mathbf{s}, \mathbf{a}, \mathbf{w})$$

probability (density) of transitioning from state \mathbf{s} to state \mathbf{s}_+ when taking action \mathbf{a}

- Stochastic policy (includes deterministic as special case)

$$\pi[\mathbf{a} | \mathbf{s}] \in \mathbb{R}_+ \quad \text{deterministic: } \mathbf{a} = \pi(\mathbf{s})$$

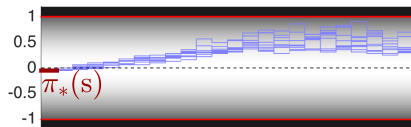
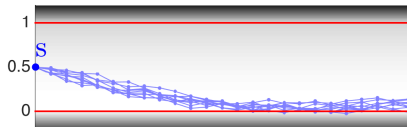
assigns the probability (density) of taking action \mathbf{a} for a given state \mathbf{s}

- Expected discounted cost (return):

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

where \mathbf{a}_k is drawn from policy π . Initial conditions \mathbf{s}_0 can be fixed or random.

- Discount factor $\gamma \in [0, 1]$



RL in a Nutshell

- Markov Chain defined by

$$\mathbb{P}[\mathbf{s}_+ | \mathbf{s}, \mathbf{a}] \quad \text{equivalent of} \quad \mathbf{s}_+ = \mathbf{f}(\mathbf{s}, \mathbf{a}, \mathbf{w})$$

probability (density) of transitioning from state \mathbf{s} to state \mathbf{s}_+ when taking action \mathbf{a}

- Stochastic policy (includes deterministic as special case)

$$\pi[\mathbf{a} | \mathbf{s}] \in \mathbb{R}_+ \quad \text{deterministic: } \mathbf{a} = \pi(\mathbf{s})$$

assigns the probability (density) of taking action \mathbf{a} for a given state \mathbf{s}

- Expected discounted cost (return):

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

where \mathbf{a}_k is drawn from policy π . Initial conditions \mathbf{s}_0 can be fixed or random.

- Discount factor $\gamma \in [0, 1]$
- Markov Decision Process (MDP): find π_* solution of

$$\pi_* := \arg \min_{\pi} J(\pi)$$

RL in a Nutshell

- Markov Chain defined by

$$\mathbb{P}[\mathbf{s}_+ | \mathbf{s}, \mathbf{a}] \quad \text{equivalent of} \quad \mathbf{s}_+ = \mathbf{f}(\mathbf{s}, \mathbf{a}, \mathbf{w})$$

probability (density) of transitioning from state \mathbf{s} to state \mathbf{s}_+ when taking action \mathbf{a}

- Stochastic policy (includes deterministic as special case)

$$\pi[\mathbf{a} | \mathbf{s}] \in \mathbb{R}_+ \quad \text{deterministic: } \mathbf{a} = \pi(\mathbf{s})$$

assigns the probability (density) of taking action \mathbf{a} for a given state \mathbf{s}

- Expected discounted cost (return):

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

where \mathbf{a}_k is drawn from policy π . Initial conditions \mathbf{s}_0 can be fixed or random.

- Discount factor $\gamma \in [0, 1]$
- Markov Decision Process (MDP): find π_* solution of

$$\pi_* := \arg \min_{\pi} J(\pi)$$

- Optimal parametrized policy π_θ given by:

$$\theta_* := \arg \min_{\theta} J(\pi_\theta)$$

RL in a Nutshell

- Markov Chain defined by

$$\mathbb{P}[\mathbf{s}_+ | \mathbf{s}, \mathbf{a}] \quad \text{equivalent of} \quad \mathbf{s}_+ = \mathbf{f}(\mathbf{s}, \mathbf{a}, \mathbf{w})$$

probability (density) of transitioning from state \mathbf{s} to state \mathbf{s}_+ when taking action \mathbf{a}

- Stochastic policy (includes deterministic as special case)

$$\pi[\mathbf{a} | \mathbf{s}] \in \mathbb{R}_+ \quad \text{deterministic: } \mathbf{a} = \pi(\mathbf{s})$$

assigns the probability (density) of taking action \mathbf{a} for a given state \mathbf{s}

- Expected discounted cost (return):

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

where \mathbf{a}_k is drawn from policy π . Initial conditions \mathbf{s}_0 can be fixed or random.

- Discount factor $\gamma \in [0, 1]$
- Markov Decision Process (MDP): find π_* solution of

$$\pi_* := \arg \min_{\pi} J(\pi)$$

- Optimal parametrized policy π_θ given by:

$$\theta_* := \arg \min_{\theta} J(\pi_\theta)$$

RL (approximately) solves the MDP in a sample-based fashion.

Value functions

- **Value function:**

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

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

Value functions

- **Value function:**

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

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

- **Note that:** $J(\pi) = \mathbb{E}_{\rho^{\mathbf{s}}} [V_{\pi}(\mathbf{s})]$, for some initial probability (density) $\rho^{\mathbf{s}}$

Value functions

- **Value function:**

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

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

- **Note that:** $J(\pi) = \mathbb{E}_{\rho^{\mathbf{s}}} [V_{\pi}(\mathbf{s})]$, for some initial probability (density) $\rho^{\mathbf{s}}$
- **Action-Value function:**

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

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

Value functions

- **Value function:**

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

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

- **Note that:** $J(\pi) = \mathbb{E}_{\rho^{\mathbf{s}}} [V_{\pi}(\mathbf{s})]$, for some initial probability (density) $\rho^{\mathbf{s}}$
- **Action-Value function:**

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

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

- **Optimal Value functions**

- ▶ Notation:

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

Value functions

- **Value function:**

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

gives the expected cost for policy π , starting from given initial conditions s

- **Note that:** $J(\pi) = \mathbb{E}_{\rho^s} [V_{\pi}(s)]$, for some initial probability (density) ρ^s
- **Action-Value function:**

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

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

- **Optimal Value functions**

- ▶ Notation:

$$V_{\star}(s) = V_{\pi_{\star}}(s)$$

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

- ▶ Properties:

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

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

Reinforcement Learning - Algorithms

Form function approximators:

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

via ad-hoc parametrization

Reinforcement Learning - Algorithms

Form function approximators:

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

via ad-hoc parametrization

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

Form function approximators:

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

via ad-hoc parametrization

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

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

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

yields policy $\pi_{\theta}(s) \approx \pi_{*}(s)$ directly

Reinforcement Learning - Algorithms

Form function approximators:

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

via ad-hoc parametrization

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

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

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

yields policy $\pi_{\theta}(s) \approx \pi_{*}(s)$ directly

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

Reinforcement Learning - Algorithms

Form function approximators:

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

via ad-hoc parametrization

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

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

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

yields policy $\pi_{\theta}(s) \approx \pi_{*}(s)$ directly

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

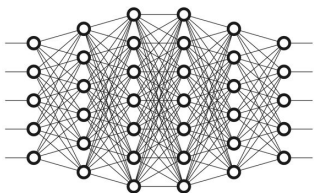
Most approaches are **derivative-based**:
we need $\nabla_{\theta} \pi_{\theta}, \nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}$

Reinforcement Learning - Algorithms

Form function approximators:

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

via ad-hoc parametrization



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

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

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

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

yields policy $\pi_{\theta}(s) \approx \pi_{*}(s)$ directly

Most approaches are **derivative-based**:
we need $\nabla_{\theta} \pi_{\theta}, \nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}$

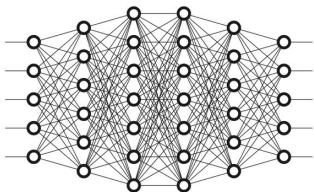
Nowadays RL typically relies on DNNs as function approximators:

Reinforcement Learning - Algorithms

Form function approximators:

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

via ad-hoc parametrization



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

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

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

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

yields policy $\pi_{\theta}(s) \approx \pi_{*}(s)$ directly

Most approaches are **derivative-based**:
we need $\nabla_{\theta} \pi_{\theta}, \nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}$

Nowadays RL typically relies on DNNs as function approximators:
difficult to understand, no strong guarantees

Form function approximators:

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

via ad-hoc parametrization

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

$$\text{s.t. } \mathbf{x}_0 = \mathbf{s},$$

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

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

$$\mathbf{h}_{\theta}^f(\mathbf{x}_N) \leq 0,$$

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

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

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

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

yields policy $\pi_{\theta}(s) \approx \pi_{*}(s)$ directly

Most approaches are **derivative-based**:
we need $\nabla_{\theta} \pi_{\theta}, \nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}$

Nowadays RL typically relies on DNNs as function approximators:
difficult to understand, no strong guarantees

Alternative: use MPC as function approximator

Form function approximators:

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

via ad-hoc parametrization

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

$$\text{s.t. } \mathbf{x}_0 = \mathbf{s},$$

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

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

$$\mathbf{h}_{\theta}^f(\mathbf{x}_N) \leq 0,$$

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

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

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

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

yields policy $\pi_{\theta}(s) \approx \pi_{*}(s)$ directly

Most approaches are **derivative-based**:
we need $\nabla_{\theta} \pi_{\theta}, \nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}$

Nowadays RL typically relies on DNNs as function approximators:
difficult to understand, no strong guarantees

Alternative: use MPC as function approximator
this provides explainability and makes it possible to guarantee safety and stability

Reinforcement Learning - Algorithms

Form function approximators:

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

via ad-hoc parametrization

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

$$\text{s.t. } \mathbf{x}_0 = \mathbf{s},$$

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

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

$$\mathbf{h}_{\theta}^f(\mathbf{x}_N) \leq 0,$$

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

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

$$\min_{\theta} \mathbb{E} \left[(Q_{*}(s, a) - Q_{\theta}(s, a))^2 \right]$$

Yields policy:

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

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

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

yields policy $\pi_{\theta}(s) \approx \pi_{*}(s)$ directly

Most approaches are **derivative-based**:
we need $\nabla_{\theta} \pi_{\theta}, \nabla_{\theta} V_{\theta}, \nabla_{\theta} Q_{\theta}$

Nowadays RL typically relies on DNNs as function approximators:
difficult to understand, no strong guarantees

Alternative: use MPC as function approximator
this provides explainability and makes it possible to guarantee safety and stability
we need to differentiate MPC!

MPC as a Function Approximator

MPC Tuning parameter θ , initial state \mathbf{s}

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

MPC as a Function Approximator

MPC Tuning parameter θ , initial state \mathbf{s}

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

- MPC delivers a parametric deterministic policy π_{θ}

MPC as a Function Approximator

MPC Tuning parameter θ , initial state \mathbf{s}

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

$$V_{\theta}^{\pi_{\theta}}(\mathbf{s}) = \min_{\mathbf{x}, \mathbf{u}} \sum_{k=0}^{N-1} \ell_{\theta}(\mathbf{x}, \mathbf{u}) + V_{\theta}^{\mathbf{f}}(\mathbf{x}_N)$$

$$\text{s.t. } \mathbf{x}_0 = \mathbf{s}$$

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

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

$$\mathbf{h}_{\theta}^{\mathbf{f}}(\mathbf{x}_N) \leq 0$$

- MPC delivers a parametric deterministic policy π_{θ}
- MPC delivers a parametric value function $V_{\theta}^{\pi_{\theta}}$

MPC as a Function Approximator

MPC Tuning parameter θ , initial state \mathbf{s} , fixed policy \mathbf{a}

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

- MPC delivers a parametric deterministic policy π_{θ}
- MPC delivers a parametric value function $V_{\theta}^{\pi_{\theta}}$
- MPC delivers a parametric action value function $Q_{\theta}^{\pi_{\theta}}$

MPC as a Function Approximator

MPC Tuning parameter θ , initial state s , fixed policy a

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

- MPC delivers a parametric deterministic policy π_{θ}
- MPC delivers a parametric value function $V_{\theta}^{\pi_{\theta}}$
- MPC delivers a parametric action value function $Q_{\theta}^{\pi_{\theta}}$
 - ▶ In practice be careful with $\mathbf{u}_0 = \mathbf{a}$: LICQ issues

MPC as a Function Approximator

MPC Tuning parameter θ , initial state s , fixed policy a

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

- MPC delivers a parametric deterministic policy π_{θ}
- MPC delivers a parametric value function $V_{\theta}^{\pi_{\theta}}$
- MPC delivers a parametric action value function $Q_{\theta}^{\pi_{\theta}}$
 - ▶ In practice be careful with $\mathbf{u}_0 = \mathbf{a}$: LICQ issues
- **All these functions can be differentiated wrt θ**
 - ▶ derivatives necessary in RL algorithms
 - ▶ need LICQ + few other conditions: mild requirements
 - ▶ if we can solve MPC we can differentiate it

MPC as a Function Approximator

MPC Tuning parameter θ , initial state s , fixed policy a

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

- MPC delivers a parametric deterministic policy π_{θ}
- MPC delivers a parametric value function $V_{\theta}^{\pi_{\theta}}$
- MPC delivers a parametric action value function $Q_{\theta}^{\pi_{\theta}}$
 - ▶ In practice be careful with $\mathbf{u}_0 = \mathbf{a}$: LICQ issues
- **All these functions can be differentiated wrt θ**
 - ▶ derivatives necessary in RL algorithms
 - ▶ need LICQ + few other conditions: mild requirements
 - ▶ if we can solve MPC we can differentiate it

MPC has all the properties required from a function approximator

- We can let RL learn the best θ
- We can let MPC enforce stability and safety guarantees in RL

MPC is a parametric NLP

MPC:

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

$$\text{s.t. } \mathbf{x}_0 = \mathbf{s}$$

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

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

$$\mathbf{h}_{\theta}^f(\mathbf{x}_N) \leq 0$$

MPC is a parametric NLP

MPC:

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

Parametric NLP:

$$\begin{aligned} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t.} \quad & \bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \\ & \bar{\mathbf{h}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \leq 0 \end{aligned}$$

where

$$\begin{aligned} \bar{\mathbf{x}} &= (\mathbf{x}, \mathbf{u}) \\ \bar{\mathbf{p}} &= \theta \end{aligned}$$

MPC is a parametric NLP

MPC:

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

Parametric NLP:

$$\begin{aligned} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t.} \quad & \bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \\ & \bar{\mathbf{h}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \leq 0 \end{aligned}$$

where

$$\begin{aligned} \bar{\mathbf{x}} &= (\mathbf{x}, \mathbf{u}) \\ \bar{\mathbf{p}} &= \theta \end{aligned}$$

Two questions require a solid answer:

- How do we solve NLPs?
- How do we differentiate NLPs?

How Do We Solve Parametric NLPs?

NLP

$$\begin{aligned} \min_{\bar{\mathbf{x}}} \quad & \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t.} \quad & \bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \\ & \bar{\mathbf{h}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \leq 0 \end{aligned}$$

Quadratic approximation

$$\begin{aligned} \min_{\Delta\bar{\mathbf{x}}} \quad & \Delta\bar{\mathbf{x}}^\top \mathbf{M} \Delta\bar{\mathbf{x}} + \mathbf{m}^\top \Delta\bar{\mathbf{x}} \\ \text{s.t.} \quad & \mathbf{G} \Delta\bar{\mathbf{x}} + \bar{\mathbf{g}} = 0 \\ & \mathbf{H} \Delta\bar{\mathbf{x}} + \bar{\mathbf{h}} \leq 0 \end{aligned}$$

Iterative procedure:

- Given $\bar{\mathbf{x}}$
- Compute quadratic approximation
- Enforce $\mathbf{M} \succ 0$
- Solve QP / linear system
- Ensure progress by computing α
- Take a step $\bar{\mathbf{x}} = \bar{\mathbf{x}} + \alpha \Delta\bar{\mathbf{x}}$

How Do We Solve Parametric NLPs?

NLP

$$\begin{aligned} \min_{\bar{\mathbf{x}}} \quad & \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t.} \quad & \bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \\ & \bar{\mathbf{h}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \leq 0 \end{aligned}$$

Iterative procedure:

- Given $\bar{\mathbf{x}}$
- Compute quadratic approximation
- Enforce $\mathbf{M} \succ 0$
- Solve QP / linear system
- Ensure progress by computing α
- Take a step $\bar{\mathbf{x}} = \bar{\mathbf{x}} + \alpha \Delta \bar{\mathbf{x}}$

Quadratic approximation

$$\begin{aligned} \min_{\Delta \bar{\mathbf{x}}} \quad & \Delta \bar{\mathbf{x}}^\top \mathbf{M} \Delta \bar{\mathbf{x}} + \mathbf{m}^\top \Delta \bar{\mathbf{x}} \\ \text{s.t.} \quad & \mathbf{G} \Delta \bar{\mathbf{x}} + \bar{\mathbf{g}} = 0 \\ & \mathbf{H} \Delta \bar{\mathbf{x}} + \bar{\mathbf{h}} \leq 0 \end{aligned}$$

Primal-dual solution

$$\bar{\mathbf{z}}^*(\bar{\mathbf{p}}) := (\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\boldsymbol{\lambda}}^*(\bar{\mathbf{p}}), \bar{\boldsymbol{\mu}}^*(\bar{\mathbf{p}}))$$

satisfies the KKT conditions

$$\mathbf{r}(\bar{\mathbf{z}}, \bar{\mathbf{p}}) := \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}, \bar{\mathbf{p}}) \\ \bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \bar{\mathbf{h}}_{\mathbb{A}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \end{bmatrix} = 0$$
$$\bar{\boldsymbol{\mu}}_{\mathbb{A}^c}^*(\bar{\mathbf{p}}) = 0$$

with active set \mathbb{A} and

$$\mathcal{L}(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\boldsymbol{\mu}}, \bar{\mathbf{p}}) := \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) + \bar{\boldsymbol{\lambda}}^\top \bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) + \bar{\boldsymbol{\mu}}^\top \bar{\mathbf{h}}(\bar{\mathbf{x}}, \bar{\mathbf{p}})$$

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ &\text{s.t. } \bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \end{aligned}$$

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned}\bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0\end{aligned}$$

KKT Conditions

$$r(\bar{\mathbf{z}}, \bar{\mathbf{p}}) := \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}(\bar{\mathbf{x}}, \bar{\boldsymbol{\lambda}}, \bar{\mathbf{p}}) \\ \bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \end{bmatrix} = 0$$

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } \bar{g}(\bar{x}, \bar{p}) &= 0 \end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

Implicit Function Theorem

Let \bar{z} be implicitly given by the (at least) \mathcal{C}^1 function

$$r(\bar{z}, \bar{p}) = 0 \quad \text{with} \quad \nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p}) \quad \text{full rank.}$$

Then \exists a \mathcal{C}^1 function $\bar{z}^*(\bar{p})$ such that

$$r(\bar{z}^*(\bar{p}), \bar{p}) = 0$$

holds in a neighbourhood of \bar{p} .

Moreover:

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } \bar{g}(\bar{x}, \bar{p}) &= 0 \end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

Implicit Function Theorem

Let \bar{z} be implicitly given by the (at least) \mathcal{C}^1 function

$$r(\bar{z}, \bar{p}) = 0 \quad \text{with} \quad \nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p}) \quad \text{full rank.}$$

Then \exists a \mathcal{C}^1 function $\bar{z}^*(\bar{p})$ such that

$$r(\bar{z}^*(\bar{p}), \bar{p}) = 0$$

holds in a neighbourhood of \bar{p} .

Moreover:

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Is this enough?

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } \bar{g}(\bar{x}, \bar{p}) &= 0 \end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

Implicit Function Theorem

Let \bar{z} be implicitly given by the (at least) \mathcal{C}^1 function

$$r(\bar{z}, \bar{p}) = 0 \quad \text{with} \quad \nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p}) \quad \text{full rank.}$$

Then \exists a \mathcal{C}^1 function $\bar{z}^*(\bar{p})$ such that

$$r(\bar{z}^*(\bar{p}), \bar{p}) = 0$$

holds in a neighbourhood of \bar{p} .

Moreover:

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Is this enough?

**Not all KKT points
are minima!**

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } \bar{g}(\bar{x}, \bar{p}) &= 0 \end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

Implicit Function Theorem

Let \bar{z} be implicitly given by the (at least) \mathcal{C}^1 function

$$r(\bar{z}, \bar{p}) = 0 \quad \text{with} \quad \nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p}) \quad \text{full rank.}$$

Then \exists a \mathcal{C}^1 function $\bar{z}^*(\bar{p})$ such that

$$r(\bar{z}^*(\bar{p}), \bar{p}) = 0$$

holds in a neighbourhood of \bar{p} .

Moreover:

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Is this enough?

**Not all KKT points
are minima!**

We need:

- LICQ
- SOSC

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } \bar{g}(\bar{x}, \bar{p}) &= 0 \end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

Implicit Function Theorem

Let \bar{z} be implicitly given by the (at least) \mathcal{C}^1 function

$$r(\bar{z}, \bar{p}) = 0 \quad \text{with} \quad \nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p}) \quad \text{full rank.}$$

Then \exists a \mathcal{C}^1 function $\bar{z}^*(\bar{p})$ such that

$$r(\bar{z}^*(\bar{p}), \bar{p}) = 0$$

holds in a neighbourhood of \bar{p} .

Moreover:

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Is this enough?

**Not all KKT points
are minima!**

We need:

- LICQ
- SOSC

Then $\nabla_{\bar{z}} r$ invertible!

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } &\bar{g}(\bar{x}, \bar{p}) = 0 \end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

How to Compute the Derivative

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Is this enough?

**Not all KKT points
are minima!**

We need:

- LICQ
- SOSC

Then $\nabla_{\bar{z}} r$ invertible!

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } &\bar{g}(\bar{x}, \bar{p}) = 0 \end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

How to Compute the Derivative

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Is this enough?

**Not all KKT points
are minima!**

We need:

- LICQ
- SOSC

Then $\nabla_{\bar{z}} r$ invertible!

Shorthand:

$$\mathcal{L}_{\bar{p}} := \mathcal{L}(\bar{x}^*(\bar{p}), \bar{\lambda}^*(\bar{p}), \bar{p})$$

$$\bar{g}_{\bar{p}} := \bar{g}(\bar{x}^*(\bar{p}), \bar{p})$$

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned} \bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } &\bar{g}(\bar{x}, \bar{p}) = 0 \end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

How to Compute the Derivative

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Observe that

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top = \begin{bmatrix} \nabla_{\bar{x}\bar{x}}^2 \mathcal{L}_{\bar{p}} & \nabla_{\bar{x}} \bar{g}_{\bar{p}} \\ \nabla_{\bar{x}} \bar{g}_{\bar{p}}^\top & 0 \end{bmatrix}$$

Is this enough?

**Not all KKT points
are minima!**

We need:

- LICQ
- SOSC

Then $\nabla_{\bar{z}} r$ invertible!

Shorthand:

$$\mathcal{L}_{\bar{p}} := \mathcal{L}(\bar{x}^*(\bar{p}), \bar{\lambda}^*(\bar{p}), \bar{p})$$

$$\bar{g}_{\bar{p}} := \bar{g}(\bar{x}^*(\bar{p}), \bar{p})$$

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned}\bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } &\bar{g}(\bar{x}, \bar{p}) = 0\end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

How to Compute the Derivative

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Observe that

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top = \begin{bmatrix} \nabla_{\bar{x}\bar{x}}^2 \mathcal{L}_{\bar{p}} & \nabla_{\bar{x}} \bar{g}_{\bar{p}} \\ \nabla_{\bar{x}} \bar{g}_{\bar{p}}^\top & 0 \end{bmatrix}$$

Then

$$\begin{bmatrix} \nabla_{\bar{x}\bar{x}}^2 \mathcal{L}_{\bar{p}} & \nabla_{\bar{x}} \bar{g}_{\bar{p}} \\ \nabla_{\bar{x}} \bar{g}_{\bar{p}}^\top & 0 \end{bmatrix} \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}_{\bar{p}} \\ \bar{g}_{\bar{p}} \end{bmatrix}$$

Is this enough?

**Not all KKT points
are minima!**

We need:

- LICQ
- SOSC

Then $\nabla_{\bar{z}} r$ invertible!

Shorthand:

$$\begin{aligned}\mathcal{L}_{\bar{p}} &:= \mathcal{L}(\bar{x}^*(\bar{p}), \bar{\lambda}^*(\bar{p}), \bar{p}) \\ \bar{g}_{\bar{p}} &:= \bar{g}(\bar{x}^*(\bar{p}), \bar{p})\end{aligned}$$

Parametric NLP Sensitivities

Parametric NLP:

$$\begin{aligned}\bar{x}^*(\bar{p}) &:= \arg \min_{\bar{x}} \bar{f}(\bar{x}, \bar{p}) \\ \text{s.t. } &\bar{g}(\bar{x}, \bar{p}) = 0\end{aligned}$$

KKT Conditions

$$r(\bar{z}, \bar{p}) := \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}(\bar{x}, \bar{\lambda}, \bar{p}) \\ \bar{g}(\bar{x}, \bar{p}) \end{bmatrix} = 0$$

How to Compute the Derivative

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} r(\bar{z}^*(\bar{p}), \bar{p})$$

Observe that

$$\nabla_{\bar{z}} r(\bar{z}^*(\bar{p}), \bar{p})^\top = \begin{bmatrix} \nabla_{\bar{x}\bar{x}}^2 \mathcal{L}_{\bar{p}} & \nabla_{\bar{x}} \bar{g}_{\bar{p}} \\ \nabla_{\bar{x}} \bar{g}_{\bar{p}}^\top & 0 \end{bmatrix}$$

Then

$$\begin{bmatrix} \nabla_{\bar{x}\bar{x}}^2 \mathcal{L}_{\bar{p}} & \nabla_{\bar{x}} \bar{g}_{\bar{p}} \\ \nabla_{\bar{x}} \bar{g}_{\bar{p}}^\top & 0 \end{bmatrix} \frac{d\bar{z}^*(\bar{p})}{d\bar{p}} = -\frac{\partial}{\partial \bar{p}} \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}_{\bar{p}} \\ \bar{g}_{\bar{p}} \end{bmatrix}$$

Compare with the (last) Newton step!

$$\begin{bmatrix} \nabla_{\bar{x}\bar{x}}^2 \mathcal{L}_{\bar{p}} & \nabla_{\bar{x}} \bar{g}_{\bar{p}} \\ \nabla_{\bar{x}} \bar{g}_{\bar{p}}^\top & 0 \end{bmatrix} \begin{bmatrix} \Delta \bar{x} \\ \Delta \bar{\lambda} \end{bmatrix} = - \begin{bmatrix} \nabla_{\bar{x}} \mathcal{L}_{\bar{p}} \\ \bar{g}_{\bar{p}} \end{bmatrix}$$

Is this enough?

**Not all KKT points
are minima!**

We need:

- LICQ
- SOSC

Then $\nabla_{\bar{z}} r$ invertible!

Shorthand:

$$\begin{aligned}\mathcal{L}_{\bar{p}} &:= \mathcal{L}(\bar{x}^*(\bar{p}), \bar{\lambda}^*(\bar{p}), \bar{p}) \\ \bar{g}_{\bar{p}} &:= \bar{g}(\bar{x}^*(\bar{p}), \bar{p})\end{aligned}$$

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}^\top} & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Important properties:

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}^\top & 0 \end{bmatrix} \frac{\mathbf{d}}{\mathbf{d}\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned}\bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0\end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}^\top} & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Important properties:

- Sensitivities are (almost) for free: KKT matrix **already factorized!**

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned}\bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0\end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}}^\top & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Important properties:

- Sensitivities are (almost) for free: KKT matrix **already factorized!**
- Sensitivity of the optimal value:

$$\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}_{\bar{\mathbf{p}}}}{\partial \bar{\mathbf{p}}}$$

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}}^\top & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Important properties:

- Sensitivities are (almost) for free: KKT matrix **already factorized!**
- Sensitivity of the optimal value:

$$\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}_{\bar{\mathbf{p}}}}{\partial \bar{\mathbf{p}}}$$

What does this mean for MPC?

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}}^\top & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Important properties:

- Sensitivities are (almost) for free: KKT matrix **already factorized!**
- Sensitivity of the optimal value:

$$\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}_{\bar{\mathbf{p}}}}{\partial \bar{\mathbf{p}}}$$

What does this mean for MPC?

- $\frac{dV_{\theta}(s)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0 \end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}}^\top & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Important properties:

- Sensitivities are (almost) for free: KKT matrix **already factorized!**
- Sensitivity of the optimal value:

$$\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}_{\bar{\mathbf{p}}}}{\partial \bar{\mathbf{p}}}$$

What does this mean for MPC?

- $\frac{dV_\theta(s)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$
- $\frac{dQ_\theta(s,a)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned}\bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0\end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}}^\top & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Important properties:

- Sensitivities are (almost) for free: KKT matrix **already factorized!**
- Sensitivity of the optimal value:

$$\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}_{\bar{\mathbf{p}}}}{\partial \bar{\mathbf{p}}}$$

What does this mean for MPC?

- $\frac{dV_\theta(s)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$
- $\frac{dQ_\theta(s,a)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$
- $\frac{d\pi_\theta(s)}{d\theta}$ from $\frac{d\bar{\mathbf{x}}^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned}\bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0\end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}}^\top & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Important properties:

- Sensitivities are (almost) for free: KKT matrix **already factorized!**
- Sensitivity of the optimal value:

$$\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}_{\bar{\mathbf{p}}}}{\partial \bar{\mathbf{p}}}$$

Sensitivities are cheap and easy to compute

What does this mean for MPC?

- $\frac{dV_\theta(s)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$
- $\frac{dQ_\theta(s,a)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$
- $\frac{d\pi_\theta(s)}{d\theta}$ from $\frac{d\bar{\mathbf{x}}^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$

Sensitivities and Implicit Functions

Parametric NLP:

$$\begin{aligned}\bar{\mathbf{x}}^*(\bar{\mathbf{p}}) &:= \arg \min_{\bar{\mathbf{x}}} \bar{f}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) \\ \text{s.t. } &\bar{\mathbf{g}}(\bar{\mathbf{x}}, \bar{\mathbf{p}}) = 0\end{aligned}$$

$$f^*(\bar{\mathbf{p}}) = \bar{f}(\bar{\mathbf{x}}^*(\bar{\mathbf{p}}), \bar{\mathbf{p}})$$

Parametric Sensitivities:

$$\begin{bmatrix} \nabla_{\bar{\mathbf{x}}\bar{\mathbf{x}}}^2 \mathcal{L}_{\bar{\mathbf{p}}} & \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}} \\ \nabla_{\bar{\mathbf{x}}\bar{\mathbf{g}}_{\bar{\mathbf{p}}}}^\top & 0 \end{bmatrix} \frac{d}{d\bar{\mathbf{p}}} \begin{bmatrix} \bar{\mathbf{x}}^*(\bar{\mathbf{p}}) \\ \bar{\lambda}^*(\bar{\mathbf{p}}) \end{bmatrix} = -\frac{\partial}{\partial \bar{\mathbf{p}}} \begin{bmatrix} \nabla_{\bar{\mathbf{x}}} \mathcal{L}_{\bar{\mathbf{p}}} \\ \bar{\mathbf{g}}_{\bar{\mathbf{p}}} \end{bmatrix}$$

Important properties:

- Sensitivities are (almost) for free: KKT matrix **already factorized!**
- Sensitivity of the optimal value:

$$\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}} = \frac{\partial \mathcal{L}_{\bar{\mathbf{p}}}}{\partial \bar{\mathbf{p}}}$$

Sensitivities are cheap and easy to compute

What does this mean for MPC?

- $\frac{dV_\theta(s)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$
- $\frac{dQ_\theta(s,a)}{d\theta}$ from $\frac{df^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$
- $\frac{d\pi_\theta(s)}{d\theta}$ from $\frac{d\bar{\mathbf{x}}^*(\bar{\mathbf{p}})}{d\bar{\mathbf{p}}}$

What about inequality constraints?

Not a real issue in practice, see parametric NLP theory

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$\hat{Q}_\theta(\mathbf{s}, \mathbf{a}) :=$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2 \\ + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0,$$

Parameter vector: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

Conditions to enforce on θ :

$\hat{Q}_\theta(\mathbf{s}, \mathbf{a}) :=$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2 \\ + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0,$$

Parameter vector: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$\hat{Q}_\theta(\mathbf{s}, \mathbf{a}) :=$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2 \\ + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\boldsymbol{\theta}) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\boldsymbol{\theta})\mathbf{x}_N + \mathbf{t}(\boldsymbol{\theta}) \leq 0,$$

Parameter vector: $\boldsymbol{\theta} = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Conditions to enforce on $\boldsymbol{\theta}$:

- Steady state

$$\mathbf{x}_r = \mathbf{A}\mathbf{x}_r + \mathbf{B}\mathbf{u}_r$$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$\hat{Q}_\theta(\mathbf{s}, \mathbf{a}) :=$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2$$

$$+ \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\boldsymbol{\theta}) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\boldsymbol{\theta})\mathbf{x}_N + \mathbf{t}(\boldsymbol{\theta}) \leq 0,$$

Parameter vector: $\boldsymbol{\theta} = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Conditions to enforce on $\boldsymbol{\theta}$:

- Steady state

$$\mathbf{x}_r = \mathbf{A}\mathbf{x}_r + \mathbf{B}\mathbf{u}_r$$

- Positive definiteness

$$\mathbf{H} = \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{R} \end{bmatrix} \succ 0$$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$\hat{Q}_\theta(\mathbf{s}, \mathbf{a}) :=$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2$$

$$+ \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\boldsymbol{\theta}) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\boldsymbol{\theta})\mathbf{x}_N + \mathbf{t}(\boldsymbol{\theta}) \leq 0,$$

Parameter vector: $\boldsymbol{\theta} = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Conditions to enforce on $\boldsymbol{\theta}$:

- Steady state

$$\mathbf{x}_r = \mathbf{A}\mathbf{x}_r + \mathbf{B}\mathbf{u}_r$$

- Positive definiteness

$$\mathbf{H} = \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{R} \end{bmatrix} \succ 0$$

- Consistent uncertainty set

$$\mathbf{M}(\mathbf{s}_{i+1} - (\mathbf{A}\mathbf{s}_i + \mathbf{B}\mathbf{a}_i + \mathbf{b})) \leq \mathbf{m}$$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$\hat{Q}_\theta(\mathbf{s}, \mathbf{a}) :=$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{array}{c} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{array} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2 + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\boldsymbol{\theta}) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\boldsymbol{\theta})\mathbf{x}_N + \mathbf{t}(\boldsymbol{\theta}) \leq 0,$$

Parameter vector: $\boldsymbol{\theta} = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Conditions to enforce on $\boldsymbol{\theta}$:

- Steady state

$$\mathbf{x}_r = \mathbf{A}\mathbf{x}_r + \mathbf{B}\mathbf{u}_r$$

- Positive definiteness

$$\mathbf{H} = \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{R} \end{bmatrix} \succ 0$$

- Consistent uncertainty set

$$\mathbf{M}(\mathbf{s}_{i+1} - (\mathbf{A}\mathbf{s}_i + \mathbf{B}\mathbf{a}_i + \mathbf{b})) \leq \mathbf{m}$$

- Terminal set includes the reference

$$\mathbf{T}(\boldsymbol{\theta})\mathbf{x}_r \leq \mathbf{t}(\boldsymbol{\theta})$$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$\hat{Q}_\theta(\mathbf{s}, \mathbf{a}) :=$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2$$

$$+ \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0,$$

Parameter vector: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Moreover:

- Riccati terminal cost and control law

$$\mathbf{K}(\theta) = (\mathbf{R} + \mathbf{B}^\top \mathbf{P}(\theta) \mathbf{B})^{-1} (\mathbf{S}^\top + \mathbf{B}^\top \mathbf{P}(\theta) \mathbf{A})$$

$$\mathbf{P}(\theta) = \mathbf{Q} + \mathbf{A}^\top \mathbf{P}(\theta) \mathbf{A} - (\mathbf{S} + \mathbf{A}^\top \mathbf{P}(\theta) \mathbf{B}) \mathbf{K}(\theta)$$

Conditions to enforce on θ :

- Steady state

$$\mathbf{x}_r = \mathbf{A}\mathbf{x}_r + \mathbf{B}\mathbf{u}_r$$

- Positive definiteness

$$\mathbf{H} = \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{R} \end{bmatrix} \succ 0$$

- Consistent uncertainty set

$$\mathbf{M}(\mathbf{s}_{i+1} - (\mathbf{A}\mathbf{s}_i + \mathbf{B}\mathbf{a}_i + \mathbf{b})) \leq \mathbf{m}$$

- Terminal set includes the reference

$$\mathbf{T}(\theta)\mathbf{x}_r \leq \mathbf{t}(\theta)$$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$\hat{Q}_\theta(s, \mathbf{a}) :=$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2 + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\boldsymbol{\theta}) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\boldsymbol{\theta})\mathbf{x}_N + \mathbf{t}(\boldsymbol{\theta}) \leq 0,$$

Parameter vector: $\boldsymbol{\theta} = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Moreover:

- Riccati terminal cost and control law

$$\mathbf{K}(\boldsymbol{\theta}) = (\mathbf{R} + \mathbf{B}^\top \mathbf{P}(\boldsymbol{\theta})\mathbf{B})^{-1}(\mathbf{S}^\top + \mathbf{B}^\top \mathbf{P}(\boldsymbol{\theta})\mathbf{A})$$

$$\mathbf{P}(\boldsymbol{\theta}) = \mathbf{Q} + \mathbf{A}^\top \mathbf{P}(\boldsymbol{\theta})\mathbf{A} - (\mathbf{S} + \mathbf{A}^\top \mathbf{P}(\boldsymbol{\theta})\mathbf{B})\mathbf{K}(\boldsymbol{\theta})$$

- Constraint tightening; RPI terminal set

$$\mathbf{f}(\boldsymbol{\theta}) = (\mathbf{c}_k(\boldsymbol{\theta}), \mathbf{T}(\boldsymbol{\theta}), \mathbf{t}(\boldsymbol{\theta}))$$

Conditions to enforce on $\boldsymbol{\theta}$:

- Steady state

$$\mathbf{x}_r = \mathbf{A}\mathbf{x}_r + \mathbf{B}\mathbf{u}_r$$

- Positive definiteness

$$\mathbf{H} = \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{R} \end{bmatrix} \succ 0$$

- Consistent uncertainty set

$$\mathbf{M}(\mathbf{s}_{i+1} - (\mathbf{A}\mathbf{s}_i + \mathbf{B}\mathbf{a}_i + \mathbf{b})) \leq \mathbf{m}$$

- Terminal set includes the reference

$$\mathbf{T}(\boldsymbol{\theta})\mathbf{x}_r \leq \mathbf{t}(\boldsymbol{\theta})$$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$$\hat{Q}_\theta(\mathbf{s}, \mathbf{a}) :=$$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2 + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0,$$

Parameter vector: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Moreover:

- Riccati terminal cost and control law

$$\mathbf{K}(\theta) = (\mathbf{R} + \mathbf{B}^\top \mathbf{P}(\theta) \mathbf{B})^{-1} (\mathbf{S}^\top + \mathbf{B}^\top \mathbf{P}(\theta) \mathbf{A})$$

$$\mathbf{P}(\theta) = \mathbf{Q} + \mathbf{A}^\top \mathbf{P}(\theta) \mathbf{A} - (\mathbf{S} + \mathbf{A}^\top \mathbf{P}(\theta) \mathbf{B}) \mathbf{K}(\theta)$$

- Constraint tightening; RPI terminal set

$$\mathbf{f}(\theta) = (\mathbf{c}_k(\theta), \mathbf{T}(\theta), \mathbf{t}(\theta))$$

RL problem

$$\min_{\theta} \psi(\theta)$$

$$\text{s.t. } \mathbf{x}_r = \mathbf{A}\mathbf{x}_r + \mathbf{B}\mathbf{u}_r$$

$$\mathbf{H} = \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{R} \end{bmatrix} \succ 0$$

$$\mathbf{M}(\mathbf{s}_{i+1} - (\mathbf{A}\mathbf{s}_i + \mathbf{B}\mathbf{a}_i + \mathbf{b})) \leq \mathbf{m}$$

$$\mathbf{T}(\theta)\mathbf{x}_r \leq \mathbf{t}(\theta)$$

Safe and Stabilizing RL Based on Linear Tube MPC

Tube-Based Robust MPC

$$\hat{Q}_\theta(s, \mathbf{a}) :=$$

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2 + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0,$$

Parameter vector: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$

Moreover:

- Riccati terminal cost and control law

$$\mathbf{K}(\theta) = (\mathbf{R} + \mathbf{B}^\top \mathbf{P}(\theta) \mathbf{B})^{-1} (\mathbf{S}^\top + \mathbf{B}^\top \mathbf{P}(\theta) \mathbf{A})$$

$$\mathbf{P}(\theta) = \mathbf{Q} + \mathbf{A}^\top \mathbf{P}(\theta) \mathbf{A} - (\mathbf{S} + \mathbf{A}^\top \mathbf{P}(\theta) \mathbf{B}) \mathbf{K}(\theta)$$

- Constraint tightening; RPI terminal set

$$\mathbf{f}(\theta) = (\mathbf{c}_k(\theta), \mathbf{T}(\theta), \mathbf{t}(\theta))$$

RL problem

$$\min_{\theta} \psi(\theta)$$

$$\text{s.t. } \mathbf{x}_r = \mathbf{A}\mathbf{x}_r + \mathbf{B}\mathbf{u}_r$$

$$\mathbf{H} = \begin{bmatrix} \mathbf{Q} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{R} \end{bmatrix} \succ 0$$

$$\mathbf{M}(\mathbf{s}_{i+1} - (\mathbf{A}\mathbf{s}_i + \mathbf{B}\mathbf{a}_i + \mathbf{b})) \leq \mathbf{m}$$

$$\mathbf{T}(\theta)\mathbf{x}_r \leq \mathbf{t}(\theta)$$

Derivative computation:

- apply chain rule:

$$\frac{d\psi(\theta)}{d\theta} = \frac{\partial\psi(\theta)}{\partial\theta} + \frac{\partial\psi(\theta)}{\partial\mathbf{f}(\theta)} \frac{d\mathbf{f}(\theta)}{d\theta}$$

- can be cumbersome but it's not rocket science

A Simple Example of Safe and Stabilizing RL

Tube MPC

$$\begin{aligned} \min_z \quad & \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{\mathbf{P}(\theta)}^2 \\ & + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l \\ \text{s.t.} \quad & \mathbf{x}_0 = \mathbf{s}, \quad \mathbf{u}_0 = \mathbf{a}, \\ & \mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1}, \\ & \mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1}, \\ & \mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0, \end{aligned}$$

Double integrator

$$\begin{aligned} \mathbf{s}_+ &= \begin{bmatrix} 1 & 0.1 \\ 0 & 1 \end{bmatrix} \mathbf{s} + \begin{bmatrix} 0.05 \\ 0.1 \end{bmatrix} \mathbf{a} + \mathbf{w} \\ \mathbf{s} &\in [-1, 1]^2 \quad \mathbf{a} \in [-1, 1] \end{aligned}$$

- Unknown noise set \mathbb{W}
- RL parameter: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$
- Reward

$$-R(\mathbf{s}, \mathbf{a}) = (p - 3)^2 + 0.01v^2 + 0.01a^2$$

A Simple Example of Safe and Stabilizing RL

Tube MPC

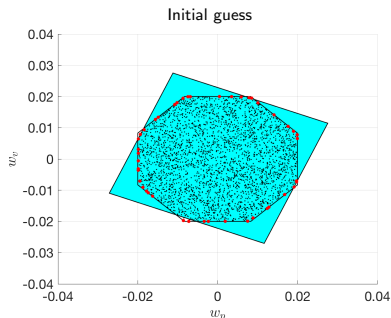
$$\begin{aligned} \min_z \quad & \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2 \\ & + \|\mathbf{x}_0\|_\Lambda^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l \\ \text{s.t.} \quad & \mathbf{x}_0 = \mathbf{s}, \quad \mathbf{u}_0 = \mathbf{a}, \\ & \mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1}, \\ & \mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1}, \\ & \mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0, \end{aligned}$$

Double integrator

$$\begin{aligned} \mathbf{s}_+ &= \begin{bmatrix} 1 & 0.1 \\ 0 & 1 \end{bmatrix} \mathbf{s} + \begin{bmatrix} 0.05 \\ 0.1 \end{bmatrix} \mathbf{a} + \mathbf{w} \\ \mathbf{s} &\in [-1, 1]^2 \quad \mathbf{a} \in [-1, 1] \end{aligned}$$

- Unknown noise set \mathbb{W}
- RL parameter: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$
- Reward

$$-R(\mathbf{s}, \mathbf{a}) = (p - 3)^2 + 0.01v^2 + 0.01a^2$$



A Simple Example of Safe and Stabilizing RL

Tube MPC

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2$$

$$+ \|\mathbf{x}_0\|_{\Lambda}^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0,$$

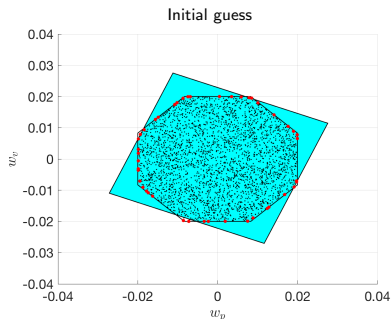
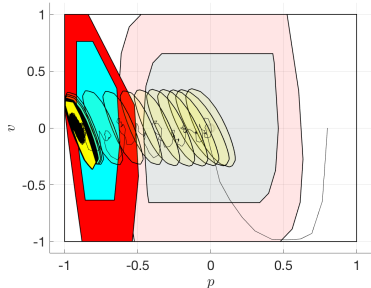
Double integrator

$$\mathbf{s}_+ = \begin{bmatrix} 1 & 0.1 \\ 0 & 1 \end{bmatrix} \mathbf{s} + \begin{bmatrix} 0.05 \\ 0.1 \end{bmatrix} \mathbf{a} + \mathbf{w}$$

$$\mathbf{s} \in [-1, 1]^2 \quad \mathbf{a} \in [-1, 1]$$

- Unknown noise set \mathbb{W}
- RL parameter: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$
- Reward

$$-R(\mathbf{s}, \mathbf{a}) = (p - 3)^2 + 0.01v^2 + 0.01a^2$$



A Simple Example of Safe and Stabilizing RL

Tube MPC

$$\min_z \sum_{k=0}^{N-1} \left\| \begin{matrix} \mathbf{x}_k - \mathbf{x}_r \\ \mathbf{u}_k - \mathbf{u}_r \end{matrix} \right\|_H^2 + \|\mathbf{x}_N - \mathbf{x}_r\|_{P(\theta)}^2$$

$$+ \|\mathbf{x}_0\|_{\Lambda}^2 + \boldsymbol{\lambda}^\top \mathbf{x}_0 + l$$

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

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \mathbf{B}\mathbf{u}_k + \mathbf{b}, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{C}\mathbf{x}_k + \mathbf{D}\mathbf{u}_k + \mathbf{c}_k(\theta) \leq 0, \quad k \in \mathbb{I}_0^{N-1},$$

$$\mathbf{T}(\theta)\mathbf{x}_N + \mathbf{t}(\theta) \leq 0,$$

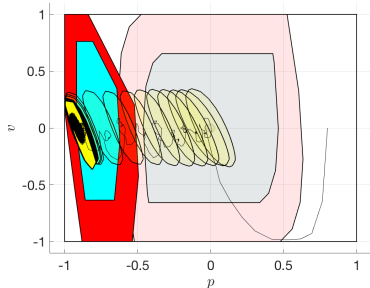
Double integrator

$$\mathbf{s}_+ = \begin{bmatrix} 1 & 0.1 \\ 0 & 1 \end{bmatrix} \mathbf{s} + \begin{bmatrix} 0.05 \\ 0.1 \end{bmatrix} \mathbf{a} + \mathbf{w}$$

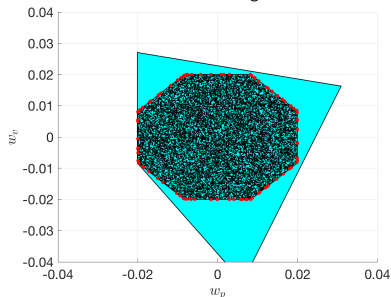
$$\mathbf{s} \in [-1, 1]^2 \quad \mathbf{a} \in [-1, 1]$$

- Unknown noise set \mathbb{W}
- RL parameter: $\theta = \{\mathbf{x}_r, \mathbf{u}_r, \mathbf{H}, \Lambda, \boldsymbol{\lambda}, l, M\}$
- Reward

$$-R(\mathbf{s}, \mathbf{a}) = (p - 3)^2 + 0.01v^2 + 0.01a^2$$



After learning



Wrap-up

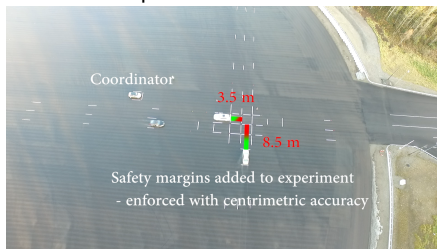
- MPC can be differentiated and that is
 - ▶ cheap
 - ▶ useful

Wrap-up

- MPC can be differentiated and that is
 - ▶ cheap
 - ▶ useful
- Differentiable MPC can be used as a function approximator within RL
 - ▶ RL can tune MPC to get the best possible performance
 - ▶ proxy to introduce safety and stability guarantees in RL
 - ▶ I have a talk on that on WeB17.2: Reinforcement Learning with Guarantees (A Tutorial on Policy Learning Methods for Advanced Controller Representations)

Wrap-up

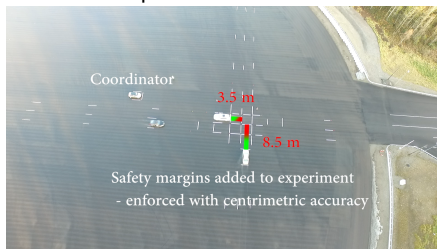
- MPC can be differentiated and that is
 - ▶ cheap
 - ▶ useful
- Differentiable MPC can be used as a function approximator within RL
 - ▶ RL can tune MPC to get the best possible performance
 - ▶ proxy to introduce safety and stability guarantees in RL
 - ▶ I have a talk on that on WeB17.2: Reinforcement Learning with Guarantees (A Tutorial on Policy Learning Methods for Advanced Controller Representations)
- Differentiable MPC can also be exploited to build distributed algorithms, e.g.:



Joint work with R. Hult, S. Gros, P. Falcone

Wrap-up

- MPC can be differentiated and that is
 - ▶ cheap
 - ▶ useful
- Differentiable MPC can be used as a function approximator within RL
 - ▶ RL can tune MPC to get the best possible performance
 - ▶ proxy to introduce safety and stability guarantees in RL
 - ▶ I have a talk on that on WeB17.2: Reinforcement Learning with Guarantees (A Tutorial on Policy Learning Methods for Advanced Controller Representations)
- Differentiable MPC can also be exploited to build distributed algorithms, e.g.:



Joint work with R. Hult, S. Gros, P. Falcone

- Parametric sensitivities are also helpful to tune tracking MPC to behave similarly to economic MPC

Our work on the topic

1. **Stability-Constrained Markov Decision Processes Using MPC**, M. Zanon, M. Palladino, S. Gros, Automatica 2022
2. **Safe Reinforcement Learning with Stability & Safety Guarantees Using Robust MPC**
Learning for MPC with Stability & Safety Guarantees, S. Gros, M. Zanon, Automatica 2022
3. **Bias Correction in Reinforcement Learning via the Deterministic Policy Gradient Method for MPC-Based Policies**, S. Gros, M. Zanon, ACC, 2021
4. **Reinforcement Learning based on MPC and the Stochastic Policy Gradient Method**, S. Gros, M. Zanon, ACC, 2021
5. **Safe Reinforcement Learning Using Robust MPC**, M. Zanon, S. Gros, IEEE TAC, 2021
6. **Safe Reinforcement Learning via Projection on a Safe Set: How to Achieve Optimality?** S. Gros, M. Zanon, A. Bemporad, IFAC World Congress 2020
7. **Reinforcement Learning for Mixed-Integer Problems Based on MPC**, S. Gros, M. Zanon, IFAC World Congress 2020
8. **Reinforcement Learning Based on Real-Time Iteration NMPC**, M. Zanon, V. Kungurtev, S. Gros, IFAC World Congress 2020
9. **Data-driven Economic NMPC using Reinforcement Learning**, S. Gros, M. Zanon, IEEE TAC, 2020
10. **Practical Reinforcement Learning of Stabilizing Economic MPC**, M. Zanon, S. Gros, A. Bemporad, European Control Conference 2019