# RL and MPC
## Safety, Stability, and some more recent results

### Sébastien Gros

Dept. of Cybernetic, NTNU
Faculty of Information Tech.

### Freiburg PhD School
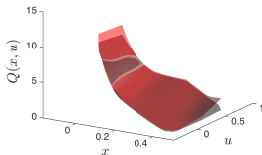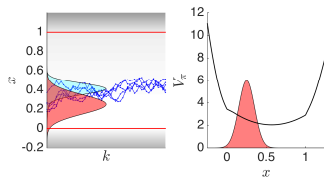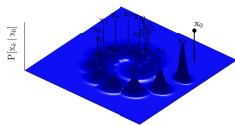
# Outline

# What are we going to discuss?

1. Learning for MPC - A focus on closed-loop performance
2. Safety & stability in Learning for MPC
3. MPC and Markov Decision Processes - When is learning beneficial?

# Outline

S. Gros (NTNU)  MPC & RL  Fall 2023  4 / 25

# Robust MPC - Uncertainty model

$$\text{True system:} \quad \mathbf{s}_+ \sim \mathbb{P}\left[\,\cdot\,|\mathbf{s}, \mathbf{a}\,\right]$$

$$\text{Deterministic model:} \quad \hat{\mathbf{s}}_+ = \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right)$$

# Robust MPC - Uncertainty model

$$\mathbb{W}_{\boldsymbol{\theta}}$$

True system: $\quad \mathbf{s}_+ \sim \mathbb{P}\left[\, \cdot \,|\mathbf{s}, \mathbf{a}\,\right]$

Deterministic model: $\quad \hat{\mathbf{s}}_+ = \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right)$



Dispersion: $\mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) + \mathbb{W}_{\boldsymbol{\theta}}$ contains the support of $\mathbb{P}\left[\, \cdot \,|\mathbf{s}, \mathbf{a}\,\right]$, i.e.

$$\mathbf{s}_+ \in \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) + \mathbb{W}_{\boldsymbol{\theta}} \tag{1}$$

with probability 1

# Robust MPC - Uncertainty model

$$\text{True system:} \quad \mathbf{s}_+ \sim \mathbb{P}\left[\,\cdot\,|\mathbf{s}, \mathbf{a}\,\right]$$

$$\text{Deterministic model:} \quad \hat{\mathbf{s}}_+ = \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right)$$



Dispersion: $\mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) + \mathbb{W}_{\boldsymbol{\theta}}$ contains the support of $\mathbb{P}\left[\,\cdot\,|\mathbf{s}, \mathbf{a}\,\right]$, i.e.

$$\mathbf{s}_+ \in \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) + \mathbb{W}_{\boldsymbol{\theta}} \qquad (1)$$

with probability 1

**Remarks**:

- Identifying $\mathbb{W}_{\boldsymbol{\theta}}$ is a set-membership identification problem, well studied

- Obviously $\mathbb{W}_{\boldsymbol{\theta}}$ is not unique

- Ensuring probability 1 from data is impossible $\rightarrow$ probabilistic guarantees

- Model parameters $\boldsymbol{\theta}$ must be such that (1) holds on every known data point

# Robust MPC - Uncertainty model

True system: $\mathbf{s}_+ \sim \mathbb{P}\left[\,\cdot\,|\mathbf{s}, \mathbf{a}\right]$

Deterministic model: $\hat{\mathbf{s}}_+ = \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right)$



Dispersion: $\mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) + \mathbb{W}_{\boldsymbol{\theta}}$ contains the support of $\mathbb{P}\left[\,\cdot\,|\mathbf{s}, \mathbf{a}\right]$, i.e.

$$\mathbf{s}_+ \in \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) + \mathbb{W}_{\boldsymbol{\theta}} \tag{1}$$

with probability 1

Condition

$$\mathbf{s}_+ - \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) \in \mathbb{W}_{\boldsymbol{\theta}}$$

for all observed triplets $(\mathbf{s}, \mathbf{a}, \mathbf{s}_+)$

$\rightarrow$ constraints on $\boldsymbol{\theta}$

**Remarks**:

- Identifying $\mathbb{W}_{\boldsymbol{\theta}}$ is a set-membership identification problem, well studied

- Obviously $\mathbb{W}_{\boldsymbol{\theta}}$ is not unique

- Ensuring probability 1 from data is impossible
  $\rightarrow$ probabilistic guarantees

- Model parameters $\boldsymbol{\theta}$ must be such that (1) holds on every known data point

# Robust MPC - Uncertainty model

True system: $\mathbf{s}_+ \sim \mathbb{P}\left[\,\cdot\,|\mathbf{s}, \mathbf{a}\right]$

Deterministic model: $\hat{\mathbf{s}}_+ = \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right)$



Dispersion: $\mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) + \mathbb{W}_{\boldsymbol{\theta}}$ contains the support of $\mathbb{P}\left[\,\cdot\,|\mathbf{s}, \mathbf{a}\right]$, i.e.

$$\mathbf{s}_+ \in \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) + \mathbb{W}_{\boldsymbol{\theta}} \qquad (1)$$

with probability 1

Condition

$$\mathbf{s}_+ - \mathbf{f}_{\boldsymbol{\theta}}\left(\mathbf{s}, \mathbf{a}\right) \in \mathbb{W}_{\boldsymbol{\theta}}$$

for all observed triplets $(\mathbf{s}, \mathbf{a}, \mathbf{s}_+)$

$\rightarrow$ constraints on $\boldsymbol{\theta}$

**Remarks**:

- Identifying $\mathbb{W}_{\boldsymbol{\theta}}$ is a set-membership identification problem, well studied

- Obviously $\mathbb{W}_{\boldsymbol{\theta}}$ is not unique

- Ensuring probability 1 from data is impossible $\rightarrow$ probabilistic guarantees

- Model parameters $\boldsymbol{\theta}$ must be such that (1) holds on every known data point
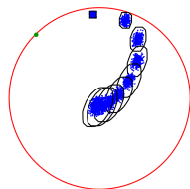
Containing the model-system mismatch becomes constraints in the parameters $\boldsymbol{\theta}$. Constraints can be readily formulated in terms of data.

## Safe policies via Robust (N)MPC

Robust (N)MPC delivers policy $\pi_\theta(x_0) = u_0^\star$ from

$$u^\star = \arg\min_u \ \max_{w \in \mathbb{W}_\theta^N} \ T_\theta(x_N) + \sum_{k=0}^{N-1} L_\theta(x_k, u_k)$$

$$\text{s.t.} \quad u_{0,\dots,N} \in \mathbb{U}$$



- $x_{0,\dots,N}$ is the propagation of the state dispersion
- max cost treats worst-case scenario, required for "classic" stability
- $w = \{w_0, \dots, w_N\}$ is the disturbance with $w_k \in \mathbb{W}_\theta$

## Safe policies via Robust (N)MPC

Robust (N)MPC delivers policy $\boldsymbol{\pi}_{\boldsymbol{\theta}}(\mathbf{x}_0) = \mathbf{u}_0^\star$ from

$$\mathbf{u}^\star = \arg\min_{\mathbf{u}} \max_{\mathbf{w} \in \mathbb{W}_{\boldsymbol{\theta}}^N} T_{\boldsymbol{\theta}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} L_{\boldsymbol{\theta}}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t.} \quad \mathbf{u}_{0,\ldots,N} \in \mathbb{U}$$

$$\mathbf{x}_{1,\ldots,N-1}(\mathbf{u}, \mathbf{x}_0, \boldsymbol{\theta}, \mathbf{w}) \in \mathbb{X}, \quad \forall \mathbf{w} \in \mathbb{W}_{\boldsymbol{\theta}}^{N-1}$$



- $\mathbf{x}_{0,\ldots,N}$ is the propagation of the state dispersion
- max cost treats worst-case scenario, required for "classic" stability
- $\mathbf{w} = \{\mathbf{w}_0, \ldots, \mathbf{w}_N\}$ is the disturbance with $\mathbf{w}_k \in \mathbb{W}_{\boldsymbol{\theta}}$
- $\mathbf{x}_{1,\ldots,N-1}(\mathbf{u}, \mathbf{x}_0, \boldsymbol{\theta}, \mathbf{w})$ are the trajectories subject to $\mathbf{w}$ and $\mathbf{f}_{\boldsymbol{\theta}}$
- $\mathbb{X}$ is the "safe" set where the state should be at all time

# Safe policies via Robust (N)MPC

Robust (N)MPC delivers policy $\pi_\theta(\mathbf{x}_0) = \mathbf{u}_0^\star$ from
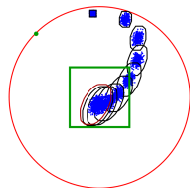
$$\mathbf{u}^\star = \arg\min_{\mathbf{u}} \max_{\mathbf{w} \in \mathbb{W}_\theta{}^N} T_\theta(\mathbf{x}_N) + \sum_{k=0}^{N-1} L_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t.} \quad \mathbf{u}_{0,\dots,N} \in \mathbb{U}$$

$$\mathbf{x}_{1,\dots,N-1}(\mathbf{u}, \mathbf{x}_0, \theta, \mathbf{w}) \in \mathbb{X}, \quad \forall \mathbf{w} \in \mathbb{W}_\theta{}^{N-1}$$

$$\mathbf{x}_N(\mathbf{u}, \mathbf{x}_0, \theta, \mathbf{w}) \in \mathbb{T}_\theta, \quad \forall \mathbf{w} \in \mathbb{W}_\theta{}^{N-1}$$

- $\mathbf{x}_{0,\dots,N}$ is the propagation of the state dispersion
- max cost treats worst-case scenario, required for "classic" stability
- $\mathbf{w} = \{\mathbf{w}_0, \dots, \mathbf{w}_N\}$ is the disturbance with $\mathbf{w}_k \in \mathbb{W}_\theta$
- $\mathbf{x}_{1,\dots,N-1}(\mathbf{u}, \mathbf{x}_0, \theta, \mathbf{w})$ are the trajectories subject to $\mathbf{w}$ and $\mathbf{f}_\theta$
- $\mathbb{X}$ is the "safe" set where the state should be at all time
- Terminal set $\mathbb{T}_\theta$ (required for recursive feasibility & stability)

# Safe policies via Robust (N)MPC

Robust (N)MPC delivers policy $\pi_\theta(\mathbf{x}_0) = \mathbf{u}_0^\star$ from
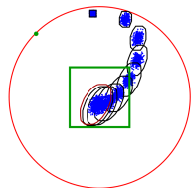
$$\mathbf{u}^\star = \arg\min_{\mathbf{u}} \ \max_{\mathbf{w} \in \mathbb{W}_\theta{}^N} \ T_\theta(\mathbf{x}_N) + \sum_{k=0}^{N-1} L_\theta(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t.} \quad \mathbf{u}_{0,\ldots,N} \in \mathbb{U}$$

$$\mathbf{x}_{1,\ldots,N-1}(\mathbf{u}, \mathbf{x}_0, \theta, \mathbf{w}) \in \mathbb{X}, \quad \forall \mathbf{w} \in \mathbb{W}_\theta{}^{N-1}$$

$$\mathbf{x}_N(\mathbf{u}, \mathbf{x}_0, \theta, \mathbf{w}) \in \mathbb{T}_\theta, \quad \forall \mathbf{w} \in \mathbb{W}_\theta{}^{N-1}$$



- $\mathbf{x}_{0,\ldots,N}$ is the propagation of the state dispersion
- max cost treats worst-case scenario, required for "classic" stability
- $\mathbf{w} = \{\mathbf{w}_0, \ldots, \mathbf{w}_N\}$ is the disturbance with $\mathbf{w}_k \in \mathbb{W}_\theta$
- $\mathbf{x}_{1,\ldots,N-1}(\mathbf{u}, \mathbf{x}_0, \theta, \mathbf{w})$ are the trajectories subject to $\mathbf{w}$ and $\mathbf{f}_\theta$
- $\mathbb{X}$ is the "safe" set where the state should be at all time
- Terminal set $\mathbb{T}_\theta$ (required for recursive feasibility & stability)
- If $\theta$ is such that $\mathbb{W}_\theta$ encloses state dispersion, **MPC yields safe policy**

# Safe policies via Robust (N)MPC

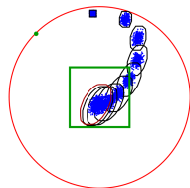Robust (N)MPC delivers policy $\pi_\theta(x_0) = u_0^\star$ from

$$u^\star = \arg\min_u \max_{w \in \mathbb{W}_\theta^N} T_\theta(x_N) + \sum_{k=0}^{N-1} L_\theta(x_k, u_k)$$

$$\text{s.t.} \quad u_{0,\ldots,N} \in \mathbb{U}$$

$$x_{1,\ldots,N-1}(u, x_0, \theta, w) \in \mathbb{X}, \quad \forall w \in \mathbb{W}_\theta^{N-1}$$

$$x_N(u, x_0, \theta, w) \in \mathbb{T}_\theta, \quad \forall w \in \mathbb{W}_\theta^{N-1}$$



- $x_{0,\ldots,N}$ is the propagation of the state dispersion
- max cost treats worst-case scenario, required for "classic" stability
- $w = \{w_0, \ldots, w_N\}$ is the disturbance with $w_k \in \mathbb{W}_\theta$
- $x_{1,\ldots,N-1}(u, x_0, \theta, w)$ are the trajectories subject to $w$ and $f_\theta$
- $\mathbb{X}$ is the "safe" set where the state should be at all time
- Terminal set $\mathbb{T}_\theta$ (required for recursive feasibility & stability)
- If $\theta$ is such that $\mathbb{W}_\theta$ encloses state dispersion, **MPC yields safe policy**

Closed-loop stability under some conditions on $\theta$ (not trivial), need $\gamma = 1$ (for now)

# Safe policies via Robust (N)MPC

Robust (N)MPC delivers policy $\boldsymbol{\pi}_{\boldsymbol{\theta}}(\mathbf{x}_0) = \mathbf{u}_0^\star$ from

$$\mathbf{u}^\star = \arg\min_{\mathbf{u}} \max_{\mathbf{w} \in \mathbb{W}_{\boldsymbol{\theta}}^N} T_{\boldsymbol{\theta}}(\mathbf{x}_N) + \sum_{k=0}^{N-1} L_{\boldsymbol{\theta}}(\mathbf{x}_k, \mathbf{u}_k)$$

$$\text{s.t. } \mathbf{u}_{0,\dots,N} \in \mathbb{U}$$

$$\mathbf{x}_{1,\dots,N-1}(\mathbf{u}, \mathbf{x}_0, \boldsymbol{\theta}, \mathbf{w}) \in \mathbb{X}, \quad \forall \mathbf{w} \in \mathbb{W}_{\boldsymbol{\theta}}^{N-1}$$

$$\mathbf{x}_N(\mathbf{u}, \mathbf{x}_0, \boldsymbol{\theta}, \mathbf{w}) \in \mathbb{T}_{\boldsymbol{\theta}}, \quad \forall \mathbf{w} \in \mathbb{W}_{\boldsymbol{\theta}}^{N-1}$$



$$\nabla_{\boldsymbol{\theta}} J = \mathbb{E}\left[\nabla_{\boldsymbol{\theta}} \boldsymbol{\pi}_{\boldsymbol{\theta}} \nabla_{\mathbf{a}} A_{\boldsymbol{\pi}_{\boldsymbol{\theta}}}\right]$$

- $\mathbf{x}_{0,\dots,N}$ is the propagation of the state dispersion
- max cost treats worst-case scenario, required for "classic" stability
- $\mathbf{w} = \{\mathbf{w}_0, \dots, \mathbf{w}_N\}$ is the disturbance with $\mathbf{w}_k \in \mathbb{W}_{\boldsymbol{\theta}}$
- $\mathbf{x}_{1,\dots,N-1}(\mathbf{u}, \mathbf{x}_0, \boldsymbol{\theta}, \mathbf{w})$ are the trajectories subject to $\mathbf{w}$ and $\mathbf{f}_{\boldsymbol{\theta}}$
- $\mathbb{X}$ is the "safe" set where the state should be at all time
- Terminal set $\mathbb{T}_{\boldsymbol{\theta}}$ (required for recursive feasibility & stability)
- If $\boldsymbol{\theta}$ is such that $\mathbb{W}_{\boldsymbol{\theta}}$ encloses state dispersion, **MPC yields safe policy**

Closed-loop stability under some conditions on $\boldsymbol{\theta}$ (not trivial), need $\gamma = 1$ (for now)

# Safe Learning via Robust MPC

**Robust NMPC parameters $\theta$**

Policy gradient

$$\nabla_\theta J = \mathbb{E}\left[\nabla_\theta \pi_\theta \nabla_u A_{\pi_\theta}\right]$$

adjusts $\theta$ for performance

Condition

$$s_+ - f(s, a, \theta) \in \mathbb{W}_\theta$$

enforces safety through $\theta$

# Safe Learning via Robust MPC

**Robust NMPC parameters $\theta$**

Policy gradient

$$\nabla_{\theta} J = \mathbb{E}\left[\nabla_{\theta}\boldsymbol{\pi}_{\theta}\nabla_{\mathbf{u}}A_{\boldsymbol{\pi}_{\theta}}\right]$$

adjusts $\theta$ for performance

Condition

$$\mathbf{s}_+ - \mathbf{f}\left(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}\right) \in \mathbb{W}_{\boldsymbol{\theta}}$$

enforces safety through $\boldsymbol{\theta}$

- No clear connection to SYSID
- Sometimes does opposite of SYSID

# Safe Learning via Robust MPC

**Robust NMPC parameters $\theta$**

**Policy gradient**

$$\nabla_{\theta} J = \mathbb{E}\left[\nabla_{\theta}\pi_{\theta}\nabla_{\mathbf{u}}A_{\pi_{\theta}}\right]$$

adjusts $\theta$ for performance

- No clear connection to SYSID
- Sometimes does opposite of SYSID

**Condition**

$$\mathbf{s}_{+} - \mathbf{f}\left(\mathbf{s}, \mathbf{a}, \theta\right) \in \mathbb{W}_{\theta}$$

enforces safety through $\theta$

- Can be interpreted as a form of SYSID (see set-membership)

# Safe Learning via Robust MPC

**Robust NMPC parameters $\theta$**

Policy gradient

$$\nabla_{\theta} J = \mathbb{E}\left[\nabla_{\theta}\pi_{\theta}\nabla_{\mathbf{u}}A_{\pi_{\theta}}\right]$$

adjusts $\theta$ for performance

- No clear connection to SYSID
- Sometimes does opposite of SYSID

Condition

$$\mathbf{s}_+ - \mathbf{f}\left(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}\right) \in \mathbb{W}_{\boldsymbol{\theta}}$$

enforces safety through $\boldsymbol{\theta}$

- Can be interpreted as a form of SYSID (see set-membership)

**How to do Safe RL?**

Classic RL steps: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha\nabla_{\boldsymbol{\theta}}J$

# Safe Learning via Robust MPC

## Robust NMPC parameters $\theta$

Policy gradient

$$\nabla_{\theta} J = \mathbb{E}\left[\nabla_{\theta} \pi_{\theta} \nabla_{\mathbf{u}} A_{\pi_{\theta}}\right]$$

adjusts $\theta$ for performance

- No clear connection to SYSID
- Sometimes does opposite of SYSID

Condition

$$\mathbf{s}_+ - \mathbf{f}\left(\mathbf{s}, \mathbf{a}, \theta\right) \in \mathbb{W}_{\theta}$$

enforces safety through $\theta$

- Can be interpreted as a form of SYSID (see set-membership)

## How to do Safe RL?

Classic RL steps: $\theta \leftarrow \theta - \alpha \nabla_{\theta} J$

Also reads as:

$$\theta \leftarrow \theta + \Delta\theta$$

$$\Delta\theta = \arg\min_{\Delta\theta} \frac{1}{2\alpha} \|\Delta\theta\|^2 + \nabla_{\theta} J^{\top} \Delta\theta$$

## Safe Learning via Robust MPC

**Robust NMPC parameters $\theta$**

Policy gradient

$$\nabla_{\boldsymbol{\theta}} J = \mathbb{E}\left[\nabla_{\boldsymbol{\theta}} \boldsymbol{\pi_\theta} \nabla_{\mathbf{u}} A_{\boldsymbol{\pi_\theta}}\right]$$

adjusts $\boldsymbol{\theta}$ for performance

- No clear connection to SYSID
- Sometimes does opposite of SYSID

Condition

$$\mathbf{s}_+ - \mathbf{f}\left(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}\right) \in \mathbb{W}_{\boldsymbol{\theta}}$$

enforces safety through $\boldsymbol{\theta}$

- Can be interpreted as a form of SYSID (see set-membership)

**How to do Safe RL?**

Classic RL steps: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} J$
Also reads as:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta\boldsymbol{\theta}$$
$$\Delta\boldsymbol{\theta} = \arg\min_{\Delta\boldsymbol{\theta}} \frac{1}{2\alpha} \|\Delta\boldsymbol{\theta}\|^2 + \nabla_{\boldsymbol{\theta}} J^\top \Delta\boldsymbol{\theta}$$

Safe RL steps $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta\boldsymbol{\theta}$:

$$\Delta\boldsymbol{\theta} = \arg\min_{\Delta\boldsymbol{\theta}} \frac{1}{2\alpha} \|\Delta\boldsymbol{\theta}\|^2 + \nabla_{\boldsymbol{\theta}} J^\top \Delta\boldsymbol{\theta}$$
$$\text{s.t. } \mathbf{s}_+ - \mathbf{f}\left(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta} + \Delta\boldsymbol{\theta}\right) \in \mathbb{W}_{\boldsymbol{\theta}+\Delta\boldsymbol{\theta}}$$
$$\forall\left(\mathbf{s}, \mathbf{a}, \mathbf{s}_+\right) \text{ in data set}$$

# Safe Learning via Robust MPC

**Robust NMPC parameters $\theta$**

Policy gradient

$$\nabla_{\boldsymbol{\theta}} J = \mathbb{E}\left[\nabla_{\boldsymbol{\theta}} \boldsymbol{\pi}_{\boldsymbol{\theta}} \nabla_{\mathbf{u}} A_{\boldsymbol{\pi}_{\boldsymbol{\theta}}}\right]$$

adjusts $\boldsymbol{\theta}$ for performance

- No clear connection to SYSID
- Sometimes does opposite of SYSID

Condition

$$\mathbf{s}_+ - \mathbf{f}\left(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta}\right) \in \mathbb{W}_{\boldsymbol{\theta}}$$

enforces safety through $\boldsymbol{\theta}$

- Can be interpreted as a form of SYSID (see set-membership)

**How to do Safe RL?**

Classic RL steps: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} J$

Also reads as:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta\boldsymbol{\theta}$$
$$\Delta\boldsymbol{\theta} = \arg\min_{\Delta\boldsymbol{\theta}} \frac{1}{2\alpha} \|\Delta\boldsymbol{\theta}\|^2 + \nabla_{\boldsymbol{\theta}} J^\top \Delta\boldsymbol{\theta}$$

Safe RL steps $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta\boldsymbol{\theta}$:

$$\Delta\boldsymbol{\theta} = \arg\min_{\Delta\boldsymbol{\theta}} \frac{1}{2\alpha} \|\Delta\boldsymbol{\theta}\|^2 + \nabla_{\boldsymbol{\theta}} J^\top \Delta\boldsymbol{\theta}$$
$$\text{s.t. } \mathbf{s}_+ - \mathbf{f}\left(\mathbf{s}, \mathbf{a}, \boldsymbol{\theta} + \Delta\boldsymbol{\theta}\right) \in \mathbb{W}_{\boldsymbol{\theta}+\Delta\boldsymbol{\theta}}$$
$$\forall \left(\mathbf{s}, \mathbf{a}, \mathbf{s}_+\right) \text{ in data set}$$

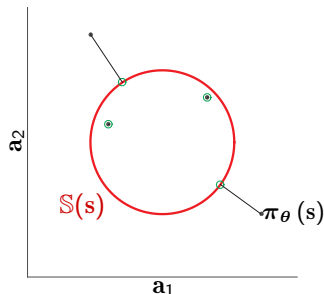**Safe RL steps seek performance under safety constraints**
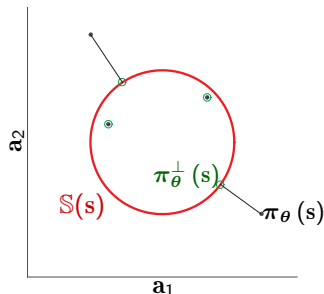
## Safety filters - Safe RL via projections

- RL can discover policy parameters $\theta$ such that policy $\pi_\theta(s)$ has good closed-loop performances, ignoring safety (e.g. $\pi_\theta$ stems from a DNN). "Learning" safety implicitly is difficult.

# Safety filters - Safe RL via projections

- RL can discover policy parameters $\theta$ such that policy $\pi_\theta(s)$ has good closed-loop performances, ignoring safety (e.g. $\pi_\theta$ stems from a DNN). "Learning" safety implicitly is difficult.

- If safe action set $\mathbb{S}(s)$ is somehow known, then we can

  ▶ follow learned policy $\pi_\theta(s)$ when

  $$\pi_\theta(s) \in \mathbb{S}(s)$$

  ▶ take "closest" action $a \in \mathbb{S}(s)$ when

  $$\pi_\theta(s) \notin \mathbb{S}(s)$$

# Safety filters - Safe RL via projections

- RL can discover policy parameters $\theta$ such that policy $\pi_\theta(s)$ has good closed-loop performances, ignoring safety (e.g. $\pi_\theta$ stems from a DNN). "Learning" safety implicitly is difficult.

- If safe action set $\mathbb{S}(s)$ is somehow known, then we can

  - follow learned policy $\pi_\theta(s)$ when
  $$\pi_\theta(s) \in \mathbb{S}(s)$$

  - take "closest" action $a \in \mathbb{S}(s)$ when
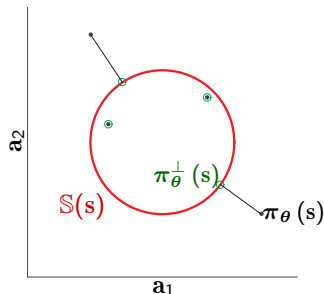  $$\pi_\theta(s) \notin \mathbb{S}(s)$$

More formally, safe policy e.g. reads as...

$$\pi_\theta^\perp(s) = \arg\min_a \quad \|a - \pi_\theta(s)\|^2$$
$$\text{s.t.} \quad a \in \mathbb{S}(s)$$

## Safety filters - Safe RL via projections

- RL can discover policy parameters $\theta$ such that policy $\pi_\theta(s)$ has good closed-loop performances, ignoring safety (e.g. $\pi_\theta$ stems from a DNN). "Learning" safety implicitly is difficult.

- If safe action set $\mathbb{S}(s)$ is somehow known, then we can

  ▶ follow learned policy $\pi_\theta(s)$ when
  $$\pi_\theta(s) \in \mathbb{S}(s)$$

  ▶ take "closest" action $a \in \mathbb{S}(s)$ when
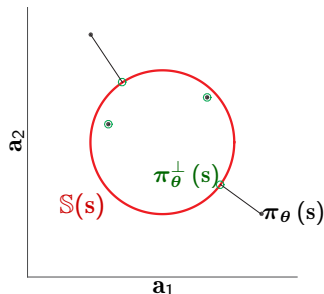  $$\pi_\theta(s) \notin \mathbb{S}(s)$$

More formally, safe policy e.g. reads as...



$$\pi_\theta^\perp(s) = \arg\min_a \quad \|a - \pi_\theta(s)\|^2$$
$$\text{s.t.} \quad a \in \mathbb{S}(s)$$

**Is that a good idea?**

# Safety filters - Safe RL via projections

- RL can discover policy parameters $\theta$ such that policy $\pi_\theta(s)$ has good closed-loop performances, ignoring safety (e.g. $\pi_\theta$ stems from a DNN). "Learning" safety implicitly is difficult.

- If safe action set $\mathbb{S}(s)$ is somehow known, then we can

  - follow learned policy $\pi_\theta(s)$ when
    $$\pi_\theta(s) \in \mathbb{S}(s)$$

  - take "closest" action $a \in \mathbb{S}(s)$ when
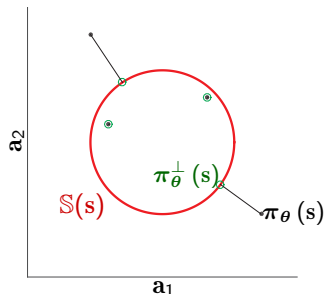    $$\pi_\theta(s) \notin \mathbb{S}(s)$$



More formally, safe policy e.g. reads as...

$$\pi_\theta^\perp(s) = \arg\min_a \quad \|a - \pi_\theta(s)\|^2$$
$$\text{s.t.} \quad a \in \mathbb{S}(s)$$

**Is that a good idea?** It depends...

# Safety filters - Safe RL via projections

- RL can discover policy parameters $\theta$ such that policy $\pi_\theta(s)$ has good closed-loop performances, ignoring safety (e.g. $\pi_\theta$ stems from a DNN). "Learning" safety implicitly is difficult.

- If safe action set $\mathbb{S}(s)$ is somehow known, then we can

  ▸ follow learned policy $\pi_\theta(s)$ when

  $$\pi_\theta(s) \in \mathbb{S}(s)$$

  ▸ take "closest" action $a \in \mathbb{S}(s)$ when

  $$\pi_\theta(s) \notin \mathbb{S}(s)$$

More formally, safe policy e.g. reads as...



$$\pi_\theta^\perp(s) = \arg\min_a \quad \|a - \pi_\theta(s)\|^2$$
$$\text{s.t.} \quad a \in \mathbb{S}(s)$$

- Built from "Robust MPC" methods?
- Interaction with learning?

**Is that a good idea?** It depends...

# Safety filters - How to obtain optimality?

**Q learning**: $Q_\theta \approx Q^\star$ learned via classic RL, ignoring safety.

## Safety filters - How to obtain optimality?

**Q learning**: $Q_\theta \approx Q^\star$ learned via classic RL, ignoring safety. Then

$$\pi_\theta^\perp(\mathbf{s}) = \arg\min_{\mathbf{a}} \quad \|\mathbf{a} - \pi_\theta(\mathbf{s})\|^2$$
$$\text{s.t.} \quad \mathbf{a} \in \mathbb{S}(\mathbf{s})$$

where

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

## Safety filters - How to obtain optimality?

**Q learning**: $Q_\theta \approx Q^\star$ learned via classic RL, ignoring safety. Then

$$\pi_\theta^\perp (\mathbf{s}) = \arg \min_{\mathbf{a}} \quad \|\mathbf{a} - \pi_\theta (\mathbf{s})\|^2$$
$$\text{s.t.} \quad \mathbf{a} \in \mathbb{S}(\mathbf{s})$$

where

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

**yields suboptimal policy $\pi_\theta^\perp$**

# Safety filters - How to obtain optimality?

**Q learning**: $Q_\theta \approx Q^\star$ learned via classic RL, ignoring safety. Then

$$\pi_\theta^\perp (\mathbf{s}) = \arg\min_\mathbf{a} \quad Q_\theta (\mathbf{s}, \mathbf{a})$$
$$\text{s.t.} \quad \mathbf{a} \in \mathbb{S}(\mathbf{s})$$

instead of a least-squares approach. Provably optimal (safe) policy.

# Safety filters - How to obtain optimality?

**Q learning**: $Q_\theta \approx Q^\star$ learned via classic RL, ignoring safety. Then

$$\pi_\theta^\perp (s) = \arg\min_a \quad Q_\theta (s, a)$$
$$\text{s.t.} \quad a \in \mathbb{S}(s)$$

instead of a least-squares approach. Provably optimal (safe) policy.

**Deterministic Policy gradient** (actor-critic): the "regular expression"

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}\left[\nabla_\theta \pi_\theta \nabla_a A_{\pi_\theta^\perp}\right]$$

yields incorrect gradients

## Safety filters - How to obtain optimality?

**Q learning**: $Q_\theta \approx Q^\star$ learned via classic RL, ignoring safety. Then

$$\boldsymbol{\pi}_\theta^\perp (\mathbf{s}) = \arg \min_{\mathbf{a}} \quad Q_\theta (\mathbf{s}, \mathbf{a})$$
$$\text{s.t.} \quad \mathbf{a} \in \mathbb{S} (\mathbf{s})$$

instead of a least-squares approach. Provably optimal (safe) policy.

**Deterministic Policy gradient** (actor-critic): make sure to evaluate the gradient using

$$\nabla_\theta J \left( \boldsymbol{\pi}_\theta^\perp \right) = \mathbb{E} \left[ \nabla_\theta \boldsymbol{\pi}_\theta^\perp \nabla_\mathbf{a} A_{\boldsymbol{\pi}_\theta^\perp} \right] \qquad \text{where} \qquad \begin{aligned} \boldsymbol{\pi}_\theta^\perp (\mathbf{s}) &= \arg \min_{\mathbf{a}} \quad \| \mathbf{a} - \boldsymbol{\pi}_\theta (\mathbf{s}) \|^2 \\ &\text{s.t.} \quad \mathbf{a} \in \mathbb{S} (\mathbf{s}) \end{aligned}$$

i.e. **account for projection (**$\Rightarrow$**differentiate NLP). Provably correct gradients**.
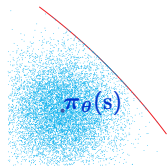
## Safety filters - How to obtain optimality?

**Q learning**: $Q_\theta \approx Q^\star$ learned via classic RL, ignoring safety. Then

$$\pi_\theta^\perp (\mathbf{s}) = \arg \min_{\mathbf{a}} \quad Q_\theta (\mathbf{s}, \mathbf{a})$$
$$\text{s.t.} \quad \mathbf{a} \in \mathbb{S} (\mathbf{s})$$

instead of a least-squares approach. Provably optimal (safe) policy.

**Deterministic Policy gradient** (actor-critic): make sure to evaluate the gradient using

$$\nabla_\theta J \left( \pi_\theta^\perp \right) = \mathbb{E} \left[ \nabla_\theta \pi_\theta^\perp \nabla_\mathbf{a} A_{\pi_\theta^\perp} \right]$$
where
$$\pi_\theta^\perp (\mathbf{s}) = \arg \min_{\mathbf{a}} \quad \| \mathbf{a} - \pi_\theta (\mathbf{s}) \|^2$$
$$\text{s.t.} \quad \mathbf{a} \in \mathbb{S} (\mathbf{s})$$

i.e. **account for projection** (⇒**differentiate NLP**). Provably **correct gradients**.

**Stochastic policy gradient**: where $\pi_\theta$ is a probability density over the actions

$$\nabla_\theta J \left( \pi_\theta^\perp \right) = \mathbb{E} \left[ \log \nabla_\theta \pi_\theta A_{\pi_\theta^\perp} \right]$$

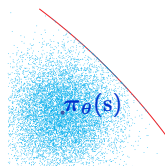i.e. **do not account for projection**. Provably **correct gradients**.

**Learning requires exploration. E.g. apply $a = \pi_\theta(s) + d$ to the real system where $d$ is a "disturbance"**

## Safe (feasible) exploration with MPC

**Learning requires exploration. E.g. apply $\mathbf{a} = \boldsymbol{\pi}_\theta(\mathbf{s}) + \mathbf{d}$ to the real system where $\mathbf{d}$ is a "disturbance"**



$\boldsymbol{\pi}_\theta(\mathbf{s})$

Explore while keeping feasibility?

- Clearly an arbitrary "policy disturbance" $\boldsymbol{\pi}_\theta(\mathbf{s}) + \mathbf{d}$ is a poor idea...
- NLP-based policy: "disturb" the cost function instead! (different options)

# Safe (feasible) exploration with MPC

**Learning requires exploration. E.g. apply $\mathbf{a} = \boldsymbol{\pi}_\theta\left(\mathbf{s}\right) + \mathbf{d}$ to the real system where $\mathbf{d}$ is a "disturbance"**



Explore while keeping feasibility?

- Clearly an arbitrary "policy disturbance" $\boldsymbol{\pi}_\theta\left(\mathbf{s}\right) + \mathbf{d}$ is a poor idea...
- NLP-based policy: "disturb" the cost function instead! (different options)

Feasible exploration: $\boldsymbol{\pi}_\theta^{\mathrm{e}}(\mathbf{s}) = \mathbf{a}_0^\star$:

$$\min_{\mathbf{x},\mathbf{u}} \quad T\left(\mathbf{x}_N\right) - \mathbf{d}^\top \mathbf{u}_0 + \sum_{k=0}^{N-1} L\left(\mathbf{x}_k, \mathbf{u}_k\right)$$

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

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

satisfies the constraints by construction

## Safe (feasible) exploration with MPC

**Learning requires exploration. E.g. apply**
$\mathbf{a} = \pi_\theta(\mathbf{s}) + \mathbf{d}$ **to the real system where** $\mathbf{d}$ **is a "disturbance"**
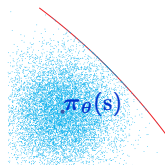

$\pi_\theta(\mathbf{s})$

Explore while keeping feasibility?

- Clearly an arbitrary "policy disturbance" $\pi_\theta(\mathbf{s}) + \mathbf{d}$ is a poor idea...
- NLP-based policy: "disturb" the cost function instead! (different options)

Feasible exploration: $\pi_\theta^{\mathrm{e}}(\mathbf{s}) = \mathbf{a}_0^\star$:

$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) - \mathbf{d}^\top \mathbf{u}_0 + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$
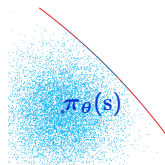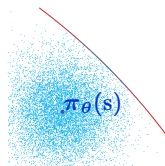
$$\mathrm{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$

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

satisfies the constraints by construction

# Safe (feasible) exploration with MPC

**Learning requires exploration. E.g. apply $\mathbf{a} = \boldsymbol{\pi}_\theta\left(\mathbf{s}\right) + \mathbf{d}$ to the real system where $\mathbf{d}$ is a "disturbance"**



Explore while keeping feasibility?

- Clearly an arbitrary "policy disturbance" $\boldsymbol{\pi}_\theta\left(\mathbf{s}\right) + \mathbf{d}$ is a poor idea...
- NLP-based policy: "disturb" the cost function instead! (different options)

Feasible exploration: $\boldsymbol{\pi}_\theta^{\mathrm{e}}(\mathbf{s}) = \mathbf{a}_0^\star$:

$$\min_{\mathbf{x}, \mathbf{u}} \quad T\left(\mathbf{x}_N\right) - \mathbf{d}^\top \mathbf{u}_0 + \sum_{k=0}^{N-1} L\left(\mathbf{x}_k, \mathbf{u}_k\right)$$

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

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

satisfies the constraints by construction

# Safe (feasible) exploration with MPC

**Learning requires exploration. E.g. apply $\mathbf{a} = \boldsymbol{\pi}_\theta(\mathbf{s}) + \mathbf{d}$ to the real system where $\mathbf{d}$ is a "disturbance"**



Explore while keeping feasibility?

- Clearly an arbitrary "policy disturbance" $\boldsymbol{\pi}_\theta(\mathbf{s}) + \mathbf{d}$ is a poor idea...
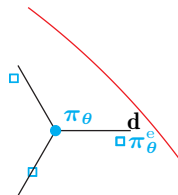- NLP-based policy: "disturb" the cost function instead! (different options)

Feasible exploration: $\boldsymbol{\pi}_\theta^{\mathrm{e}}(\mathbf{s}) = \mathbf{a}_0^\star$:

$$\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_N) - \mathbf{d}^\top \mathbf{u}_0 + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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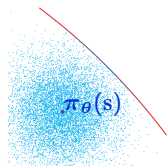
$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

satisfies the constraints by construction

# Safe (feasible) exploration with MPC

**Learning requires exploration. E.g. apply**
$\mathbf{a} = \boldsymbol{\pi}_\theta(\mathbf{s}) + \mathbf{d}$ **to the real system where** $\mathbf{d}$ **is a "disturbance"**



Explore while keeping feasibility?

- Clearly an arbitrary "policy disturbance" $\boldsymbol{\pi}_\theta(\mathbf{s}) + \mathbf{d}$ is a poor idea...
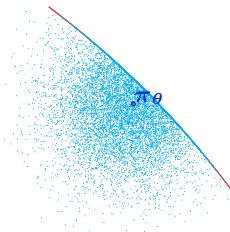- NLP-based policy: "disturb" the cost function instead! (different options)

Feasible exploration: $\boldsymbol{\pi}_\theta^{\mathrm{e}}(\mathbf{s}) = \mathbf{a}_0^\star$:

$$\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_N) - \mathbf{d}^\top \mathbf{u}_0 + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

satisfies the constraints by construction

# Safe (feasible) exploration with MPC

**Learning requires exploration. E.g. apply**
$\mathbf{a} = \boldsymbol{\pi}_\theta(\mathbf{s}) + \mathbf{d}$ **to the real system where** $\mathbf{d}$ **is a "disturbance"**



$\boldsymbol{\pi}_\theta(\mathbf{s})$

Explore while keeping feasibility?

- Clearly an arbitrary "policy disturbance" $\boldsymbol{\pi}_\theta(\mathbf{s}) + \mathbf{d}$ is a poor idea...
- NLP-based policy: "disturb" the cost function instead! (different options)

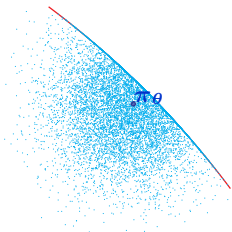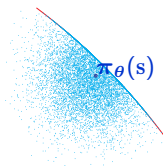Feasible exploration: $\boldsymbol{\pi}_\theta^{\mathrm{e}}(\mathbf{s}) = \mathbf{a}_0^\star$:

$$\min_{\mathbf{x}, \mathbf{u}} \quad T(\mathbf{x}_N) - \mathbf{d}^\top \mathbf{u}_0 + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$

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

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

satisfies the constraints by construction

**Remarks**:

- Exploration $\mathbf{e} = \boldsymbol{\pi}_\theta^{\mathrm{e}} - \boldsymbol{\pi}_\theta$ is not centred-isotropic
- Can create some technical issues with actor-critic methods (linear compatible $A_{\boldsymbol{\pi}_\theta}$), yields biased policy gradient estimation
- Bias seems small in practice

# Outline

## Stability of MPC

**Policy $\pi^{\mathrm{MPC}}$ from**

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

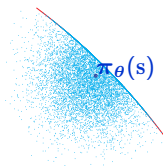$$\text{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k), \ \mathbf{x}_0 = \mathbf{s}$$

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

# Stability of MPC

**Policy $\pi^{\mathrm{MPC}}$ from**

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

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

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

MPC scheme is (nominally) stabilizing if there is $\lambda$ such that

$$\ell\left(\mathbf{s},\mathbf{a}\right) := L\left(\mathbf{s},\mathbf{a}\right) + \lambda\left(\mathbf{s}\right) - \lambda\left(\mathbf{f}\left(\mathbf{s},\mathbf{a}\right)\right) \geq \kappa\left(\|\mathbf{s} - \mathbf{s}_{\mathrm{s}}\|\right), \quad \forall\, \mathbf{s},\mathbf{a}$$

where $\kappa$ is $\mathrm{K}_\infty$ (+conditions on $T$)

## Stability of MPC

**Policy $\pi^{\mathrm{MPC}}$ from**

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

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

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

**Equivalent MPC**

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

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

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

brings us back to classic stability theory

MPC scheme is (nominally) stabilizing if there is $\lambda$ such that

$$\ell(\mathbf{s},\mathbf{a}) := L(\mathbf{s},\mathbf{a}) + \lambda(\mathbf{s}) - \lambda(\mathbf{f}(\mathbf{s},\mathbf{a})) \geq \kappa(\|\mathbf{s} - \mathbf{s}_\mathrm{s}\|), \quad \forall \, \mathbf{s}, \mathbf{a}$$

where $\kappa$ is $\mathrm{K}_\infty$ (+conditions on $T$)

## Stability of MPC

**Policy $\pi^{\mathrm{MPC}}$ from**

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

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

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

**Equivalent MPC**

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

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

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

brings us back to classic stability theory

MPC scheme is (nominally) stabilizing if there is $\lambda$ such that

$$\ell(\mathbf{s}, \mathbf{a}) := L(\mathbf{s}, \mathbf{a}) + \lambda(\mathbf{s}) - \lambda(\mathbf{f}(\mathbf{s}, \mathbf{a})) \geq \kappa(\|\mathbf{s} - \mathbf{s}_{\mathrm{s}}\|), \quad \forall \mathbf{s}, \mathbf{a}$$

where $\kappa$ is $\mathrm{K}_\infty$ (+conditions on $T$)

**Remarks**

- Modifying the MPC cost is a concept already present in dissipativity theory!
- Aligned with modifying the cost for MPC performance
- $\rightarrow$ Merge the RL & stability modifications for "Stability by design"

# Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a
**stable policy** $\boldsymbol{\pi}_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\boldsymbol{\pi}_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L\left(\mathbf{s}_k, \mathbf{a}_k\right)$$

## Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)$$

**Parametrized policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ from:

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

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

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

## Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)$$

**Parametrized policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ from:

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

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

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

- Learning based on $L$
- Impose constraint:

    $$L_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a}) \geq \kappa(\|\mathbf{s} - \mathbf{s}_{\mathrm{s}}\|), \quad \forall \mathbf{s}, \mathbf{a}$$

    throughout the learning

- $L_{\boldsymbol{\theta}}$ different than $L$ from constraint & model error

# Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\boldsymbol{\pi}_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\boldsymbol{\pi}_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)$$

- Learning based on $L$
- Impose constraint:

$$L_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a}) \geq \kappa(\|\mathbf{s} - \mathbf{s}_{\mathrm{s}}\|), \quad \forall \, \mathbf{s}, \mathbf{a}$$

throughout the learning

- $L_{\boldsymbol{\theta}}$ different than $L$ from constraint & model error

**Parametrized policy** $\boldsymbol{\pi}_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ from:

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

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

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

**Theorem**: under some conditions

- $\boldsymbol{\pi}_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to \boldsymbol{\pi}_{\star}$ if $\boldsymbol{\pi}_{\star}$ is stabilizing
- $\boldsymbol{\pi}_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to$ best stabilizing policy otherwise

## Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\pi_{\theta}^{\mathrm{MPC}}$ minimizing:

$$J\left(\pi_{\theta}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L\left(\mathbf{s}_k, \mathbf{a}_k\right)$$

- Learning based on $L$
- Impose constraint:

$$L_{\theta}\left(\mathbf{s}, \mathbf{a}\right) \geq \kappa\left(\|\mathbf{s} - \mathbf{s}_{\mathrm{s}}\|\right), \quad \forall \mathbf{s}, \mathbf{a}$$

throughout the learning

- $L_{\theta}$ different than $L$ from constraint & model error

**Parametrized policy** $\pi_{\theta}^{\mathrm{MPC}}$ from:

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

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

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

**Theorem**: under some conditions

- $\pi_{\theta}^{\mathrm{MPC}} \to \pi_{\star}$ if $\pi_{\star}$ is stabilizing
- $\pi_{\theta}^{\mathrm{MPC}} \to$ best stabilizing policy otherwise

Change of philosophy from "classic" dissipativity framework:
**stability analysis $\to$ stable design**

# Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)$$

---

**Constraint**

$$L_{\boldsymbol{\theta}}(\mathbf{s}, \mathbf{a}) \geq \kappa\left(\|\mathbf{s} - \mathbf{s}_s\|\right), \quad \forall \mathbf{s}, \mathbf{a}$$

is semi-infinite programming, not trivial

**Some solutions**:

- Sum-of-Squares (SOS) prog.
- Convex $L_{\boldsymbol{\theta}}$ (+ radially unbounded)
- Something else?

---

**Parametrized policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ from:

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

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

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

---

**Theorem**: under some conditions

- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to \pi_{\star}$ if $\pi_{\star}$ is stabilizing
- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to$ best stabilizing policy otherwise

Change of philosophy from "classic" dissipativity framework:
**stability analysis $\to$ stable design**

## Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)$$

Note that $\lambda_{\boldsymbol{\theta}}$ is redundant for policy gradient, needed for Q-learning... Combining both is meaningful!

**Parametrized policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ from:

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

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

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

**Theorem**: under some conditions

- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to \pi_{\star}$ if $\pi_{\star}$ is stabilizing
- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to$ best stabilizing policy otherwise

Change of philosophy from "classic" dissipativity framework:
**stability analysis $\to$ stable design**

## Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)$$

**Extension to stable policy for MDPs?**

- Need stability with discount
- Need "stochastic dissipativity"

**Parametrized policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ from:

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

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

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

**Theorem**: under some conditions

- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to \pi_{\star}$ if $\pi_{\star}$ is stabilizing
- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to$ best stabilizing policy otherwise

Change of philosophy from "classic" dissipativity framework:
**stability analysis $\to$ stable design**

## Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)$$

**Extension to stable policy for MDPs?**

- Need stability with discount
- Need "stochastic dissipativity"

**MDP dissipativity**: (2x Automatica '22)

- Use Strong Discounted Strict Dissipativity conditions
- Form the dissipativity equations in the measure space of the MDP

**Parametrized policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ from:

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

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

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

**Theorem**: under some conditions

- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to \pi_{\star}$ if $\pi_{\star}$ is stabilizing
- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to$ best stabilizing policy otherwise

Change of philosophy from "classic" dissipativity framework:
**stability analysis $\to$ stable design**

# Stability-constrained Learning-based MPC - Deterministic case

Given arbitrary stage cost $L(\mathbf{s}, \mathbf{a})$, build a **stable policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ minimizing:

$$J\left(\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}\right) = \sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)$$

**Extension to stable policy for MDPs?**

- Need stability with discount
- Need "stochastic dissipativity"

**MDP dissipativity**: (2x Automatica '22)

- Use Strong Discounted Strict Dissipativity conditions
- Form the dissipativity equations in the measure space of the MDP

**Parametrized policy** $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}}$ from:

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

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

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

**Theorem**: under some conditions

- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to \pi_{\star}$ if $\pi_{\star}$ is stabilizing
- $\pi_{\boldsymbol{\theta}}^{\mathrm{MPC}} \to$ best stabilizing policy otherwise

Change of philosophy from "classic" dissipativity framework:
**stability analysis → stable design**

We have the maths to treat this, not yet the algorithms...

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**

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

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

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

**MDP**:

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

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

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

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**

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

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

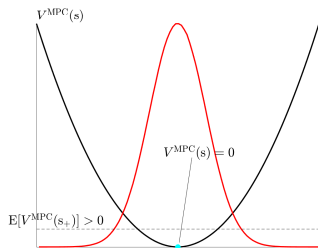$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

**MDP**:

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

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

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

**Classic stability via Lyapunov:**

- $V^{\mathrm{MPC}}(\mathbf{s})$ **decrease** along the system trajectories, i.e.

$$V^{\mathrm{MPC}}\left(\mathbf{f}\left(\mathbf{s}, \boldsymbol{\pi}^{\mathrm{MPC}}(\mathbf{s})\right)\right) < V^{\mathrm{MPC}}(\mathbf{s})$$

is ensured by construction

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**

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

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

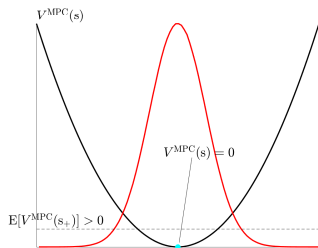$$\mathbf{h}\left(\mathbf{x}_k, \mathbf{u}_k\right) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

**MDP**:

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

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

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

**Classic stability via Lyapunov:**

- $V^{\mathrm{MPC}}\left(\mathbf{s}\right)$ **decrease** along the system trajectories, i.e.

$$V^{\mathrm{MPC}}\left(\mathbf{f}\left(\mathbf{s}, \boldsymbol{\pi}^{\mathrm{MPC}}\left(\mathbf{s}\right)\right)\right) < V^{\mathrm{MPC}}\left(\mathbf{s}\right)$$

  is ensured by construction

- What if $\mathbf{s}_+ \sim \mathbb{P}[\,.\mid \mathbf{s}, \boldsymbol{\pi}^\star\left(\mathbf{s}\right)]$ is stochastic (with know density)?

$$V^{\mathrm{MPC}}\left(\mathbf{s}_+\right) < V^{\mathrm{MPC}}\left(\mathbf{s}\right), \quad \forall \mathbf{s}$$

  in *some sense*? **Not really**... (unless strong assumptions)

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**

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

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

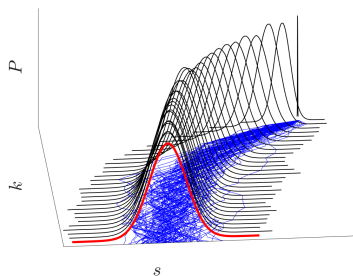$$\mathbf{h}\left(\mathbf{x}_k, \mathbf{u}_k\right) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

**MDP:**

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

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

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

**E.g. thought experiment:** $V^{\mathrm{MPC}}$ **convex, s at the minimum...**

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**

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

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

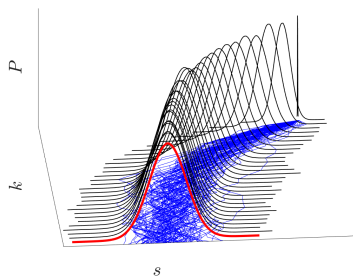$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

**MDP**:

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

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

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

**E.g. thought experiment:** $V^{\mathrm{MPC}}$ convex, $\mathbf{s}$ at the minimum...



**A Lyapunov stability theory for MDP in terms of state (beyond "stability to a set") is in general not possible.**

Yet MDPs can be stable

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**

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

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

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

**MDP**:

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

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

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

**E.g. thought experiment:** $V^{\mathrm{MPC}}$ convex, s at the minimum...



**A Lyapunov stability theory for MDP in terms of state (beyond "stability to a set") is in general not possible.**

Yet MDPs can be stable

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**

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

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

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

**MDP:**

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

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

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

**E.g. thought experiment:** $V^{\mathrm{MPC}}$ convex, s at the minimum...



**A Lyapunov stability theory for MDP in terms of state (beyond "stability to a set") is in general not possible.**

Yet MDPs can be stable

**Key idea:** Lyapunov stability in the **state measure** rather than state space

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**

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

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

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

**MDP**:

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

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

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

**Key idea**: Lyapunov stability in the **state measure** rather than state space

**Functional dissipativity**: if there is a functional $\lambda$ such that:

$$\mathcal{L}\left[\rho, \boldsymbol{\pi}\right] - \lambda\left[\rho_+\right] + \lambda\left[\rho\right] \geq \kappa\left(D\left(\rho\,||\,\rho^{\mathrm{s}}\right)\right), \qquad \mathbf{s} \sim \rho, \ \mathbf{s}_+ \sim \rho_+$$

then the state distribution $\rho$ converges to $\rho^{\mathrm{s}}$

where

- $\mathcal{L}$ is the problem cost functional, e.g. $\mathcal{L} = \mathbb{E}\left[L\left(\mathbf{s}, \mathbf{a}\right)\right]$
- $D\left(\,\cdot\,||\,\cdot\,\right)$ is a dissimilarity measure, e.g. Kullback-Liebler Divergence
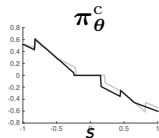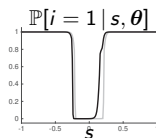- Choice of dissimilarity measure defines the form of stability

## Stability of MPC - Stochastic dynamics

**Policy $\pi_{\mathrm{MPC}}$ from**
$$\min_{\mathbf{x},\mathbf{u}} \quad T(\mathbf{x}_N) + \sum_{k=0}^{N-1} L(\mathbf{x}_k, \mathbf{u}_k)$$
$$\text{s.t.} \quad \mathbf{x}_{k+1} = \mathbf{f}(\mathbf{x}_k, \mathbf{u}_k)$$
$$\mathbf{h}(\mathbf{x}_k, \mathbf{u}_k) \leq 0, \quad \mathbf{x}_0 = \mathbf{s}$$

**MDP**:
$$\min_{\boldsymbol{\pi}} \quad \mathbb{E}_{\boldsymbol{\pi}}\left[\sum_{k=0}^{\infty} L(\mathbf{s}_k, \mathbf{a}_k)\right]$$
where $\mathbf{a}_k = \boldsymbol{\pi}(\mathbf{s}_k)$ and system dynamics
$$\mathbf{s}_{k+1} \sim \mathbb{P}\left[\,\cdot\,|\,\mathbf{s}_k, \mathbf{a}_k\,\right]$$

**Key idea**: Lyapunov stability in the **state measure** rather than state space

**Functional dissipativity**: if there is a functional $\lambda$ such that:
$$\mathcal{L}[\rho, \boldsymbol{\pi}] - \lambda[\rho_+] + \lambda[\rho] \geq \kappa\left(D\left(\rho \,||\, \rho^{\mathrm{s}}\right)\right), \qquad \mathbf{s} \sim \rho, \ \mathbf{s}_+ \sim \rho_+$$

then the state distribution $\rho$ converges to $\rho^{\mathrm{s}}$

where

- $\mathcal{L}$ is the problem cost functional, e.g. $\mathcal{L} = \mathbb{E}\left[L(\mathbf{s}, \mathbf{a})\right]$
- $D\left(\cdot \,||\, \cdot\right)$ is a dissimilarity measure, e.g. Kullback-Liebler Divergence
- Choice of dissimilarity measure defines the form of stability

**Not obvious how to use it in RL yet...**

# Outline

# RL & Mixed integer problem in MPC

**Mixed-integer problems are common. Can we do RL over Mixed-integer MPC schemes?**

*Assume mixed-integer actions*

# RL & Mixed integer problem in MPC



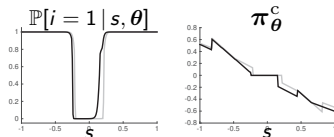Mixed-integer problems are common. Can we do RL over Mixed-integer MPC schemes?

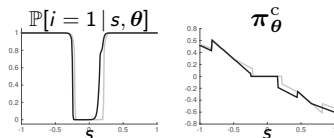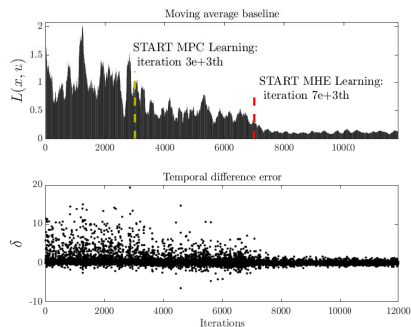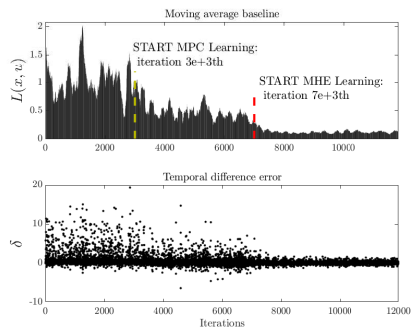*Assume mixed-integer actions*

- With Q-learning, fairly trivial... incorrect if no exploration, though

# RL & Mixed integer problem in MPC

Mixed-integer problems are common. Can we do
RL over Mixed-integer MPC schemes?



$\mathbb{P}[i = 1 \,|\, s, \theta]$      $\pi_\theta^{\mathrm{c}}$

*Assume mixed-integer actions*

- With Q-learning, fairly trivial... incorrect if no exploration, though

- For policy gradient, devil is in the details

  - ✓ Integer inputs are best treated via stochastic policy gradient
  - ✓ Continuous inputs are "best treated" via deterministic policy gradient (in the presence of constraints)
  - ✓ Propose a hybrid policy gradient method combining deterministic and stochastic policies, with corresponding compatible linear $A_{\pi_\theta}$ approximations
  - ✓ Works well on mixed-integer MPC examples

# RL & Mixed integer problem in MPC

**Mixed-integer problems are common. Can we do RL over Mixed-integer MPC schemes?**



$\mathbb{P}[i = 1 \,|\, s, \theta]$

$\pi_{\theta}^{\mathrm{c}}$

*Assume mixed-integer actions*

- With Q-learning, fairly trivial... incorrect if no exploration, though

- For policy gradient, devil is in the details

  - ✓ Integer inputs are best treated via stochastic policy gradient
  - ✓ Continuous inputs are "best treated" via deterministic policy gradient (in the presence of constraints)
  - ✓ Propose a hybrid policy gradient method combining deterministic and stochastic policies, with corresponding compatible linear $A_{\pi_{\theta}}$ approximations
  - ✓ Works well on mixed-integer MPC examples

More to be done on efficiency & control of the integer exploration

# RL & MHE-MPC



The full state of the system is often not available, or not even modelled, use observer (e.g. MHE). Can we still do RL and how?

# RL & MHE-MPC



Moving average baseline

START MPC Learning:
iteration 3e+3th

START MHE Learning:
iteration 7e+3th

Temporal difference error

The full state of the system is often not available, or not even modelled, use observer (e.g. MHE). Can we still do RL and how?

- Problem becomes POMDP when MPC model does not include all states

# RL & MHE-MPC



The full state of the system is often not available, or not even modelled, use observer (e.g. MHE). Can we still do RL and how?

- Problem becomes POMDP when MPC model does not include all states

- MHE becomes a component of the policy, must be treated in RL as well

  - ✓ RL can tune MHE and MPC jointly for closed loop performance in the context of Q learning
  - ✓ Algorithmic is simple, performances on example are good
  - ✓ The MHE tuning has a strong impact on performance (on our examples)
  - ✓ Extension to policy gradient is simple, yet to publish
  - ✓ Works also if MPC model omits some of the real states

## Tuning of the MPC "meta"-parameters

**MPC "meta"-parameters:**

- Horizon length $N$
- When to recompute control sequence (event-based MPC)

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

yields $\boldsymbol{\pi}_{\mathrm{MPC}}(\mathbf{s}_0) = \mathbf{u}_0^\star$

**Event-triggered:**

- apply input profile $\mathbf{u}_{0,\ldots,n}^\star$ until re-computation is triggered
- often used to reduce computational demand, energy, communication, etc.
- Triggering is state-based, to be tuned

# Tuning of the MPC "meta"-parameters

**MPC "meta"-parameters:**

- Horizon length $N$
- When to recompute control sequence (event-based MPC)

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

yields $\boldsymbol{\pi}_{\mathrm{MPC}}(\mathbf{s}_0) = \mathbf{u}_0^\star$

**Event-triggered:**

- apply input profile $\mathbf{u}_{0,\ldots,n}^\star$ until re-computation is triggered
- often used to reduce computational demand, energy, communication, etc.
- Triggering is state-based, to be tuned

Fairly simple idea, requires some care to be treated correctly:

- ✓ Define augmented state to preserve Markov property (essential for RL methods)
- ✓ Stochastic policy gradient methods required, must define the densities very carefully

# RL to evaluate the storage function

**Policy** $\pi_{\mathrm{MPC}}$ **from**

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

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

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

If for some $\lambda$ function:

$$L(\mathbf{s}, \mathbf{a}) + \lambda(\mathbf{s}) - \lambda(\mathbf{f}(\mathbf{s}, \mathbf{a})) \geq \kappa(\|\mathbf{s} - \mathbf{s}_{\mathrm{s}}\|), \quad \forall \, \mathbf{s}, \mathbf{a}$$

holds, then MPC scheme is stabilizing

How to evaluate $\lambda$?

- Approximate $\mathbf{f}$ as a polynomial, then Sum-of-Squares technique can be used
- We propose: parametrize $\lambda$ and evaluate it via Q-learning
- On some examples, provides a more accurate $\lambda$ than SOS
- Combination would arguably be good, to be done

# MPC Beyond State Space

**Systems with**

- $\sim$Linear dynamics
- Input-output data
- Significant stochasticity
- Modelling is difficult

# MPC Beyond State Space

**Systems with**

- $\sim$Linear dynamics
- Input-output data
- Significant stochasticity
- Modelling is difficult

**Multi-step linear predictors**

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

- Recent input-output sequence $\mathbf{u}, \mathbf{y}$
- Planned input sequence $\mathbf{u}$
- Predicted output sequence $\hat{\mathbf{y}}$

# MPC Beyond State Space

## Multi-step linear predictors

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

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

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

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

$$\mathbf{h}\left(\hat{\mathbf{y}}_k, \mathbf{u}_k\right) \leq 0$$

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

- Recent input-output sequence $\mathbf{u}, \mathbf{y}$
- Planned input sequence $\mathbf{u}$
- Predicted output sequence $\hat{\mathbf{y}}$

# MPC Beyond State Space

**Multi-step linear predictors**

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

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

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

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

$$\mathbf{h}\left(\hat{\mathbf{y}}_k, \mathbf{u}_k\right) \leq 0$$

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

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

Where $\Phi$ can be **built from past data** $\mathcal{D}$, e.g.

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

$$\text{s.t.} \quad \Phi \text{ causal}$$

or other regressions

# MPC Beyond State Space

**Multi-step linear predictors**

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

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

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

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

$$\mathbf{h}\left(\hat{\mathbf{y}}_k, \mathbf{u}_k\right) \leq 0$$

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

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

Where $\Phi$ can be **built from past data** $\mathcal{D}$, e.g.

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

$$\text{s.t.} \quad \Phi \text{ causal}$$

or other regressions

Suffers from the same
limitations as classic MPC

# MPC Beyond State Space

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

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

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

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

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

**Can we do RL? Yes!**

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

# MPC Beyond State Space

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

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

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

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

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

**Can we do RL? Yes!**

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

# MPC Beyond State Space

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

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

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

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

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

**Can we do RL? Yes!**

- RL-MPC theory applies with some twists
- State becomes $\mathbf{u}$, $\mathbf{y}$ (window of input-output)
- Modifications in principle not localized in time
- High-dimensional parameter space for RL
- Better behaved for learning than one-step simulation models *(?)*

# MPC Beyond State Space

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

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

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

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

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

**Can we do RL? Yes!**

- RL-MPC theory applies with some twists
- State becomes $\mathbf{u}$, $\mathbf{y}$ (window of input-output)
- Modifications in principle not localized in time
- High-dimensional parameter space for RL
- Better behaved for learning than one-step simulation models *(?)*

**Nonlinear extension possible. Best way to do it is to be investigated.**

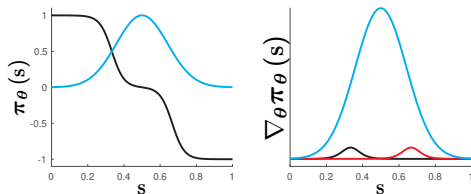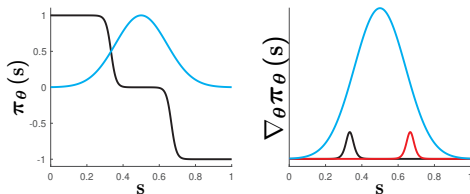# RL & MPC for "strongly economic" problems

Some policies are dominated by "switches", difficult to treat in RL because $\nabla_\theta \pi_\theta = 0$ on most of the state space. Hence

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

is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because**
$$\nabla_\theta \pi_\theta = 0 \text{ on most of the state space. Hence}$$

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere
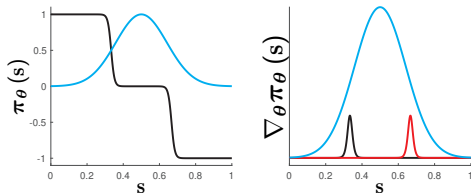
✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because $\nabla_\theta \pi_\theta = 0$ on most of the state space. Hence**

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere
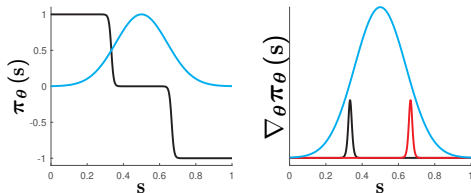
✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because**
$$\nabla_\theta \pi_\theta = 0 \text{ on most of the state space. Hence}$$

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere
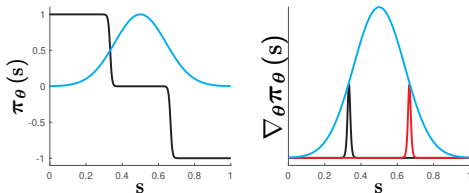
✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

> **Some policies are dominated by "switches", difficult to treat in RL because $\nabla_\theta \pi_\theta = 0$ on most of the state space. Hence**
>
> $$\nabla_\theta J(\pi_\theta) = \mathbb{E}\left[\nabla_\theta \pi_\theta \nabla_a A_{\pi_\theta}\right]$$
>
> **is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere

✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because**
$$\nabla_\theta \pi_\theta = 0 \text{ on most of the state space. Hence}$$

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere
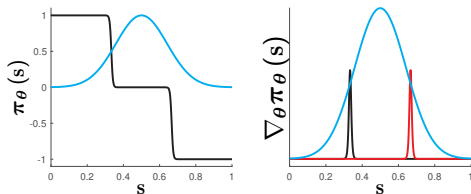
✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because**
$$\nabla_\theta \pi_\theta = 0 \text{ on most of the state space. Hence}$$

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

- ✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere
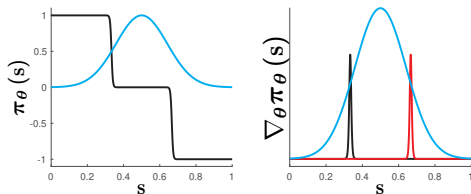- ✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because**
$$\nabla_\theta \pi_\theta = 0 \text{ on most of the state space. Hence}$$

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere
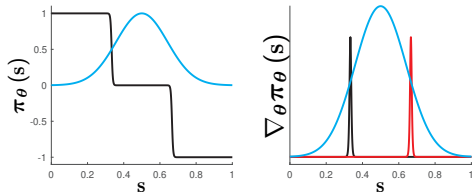
✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because**
$$\nabla_\theta \pi_\theta = 0 \text{ on most of the state space. Hence}$$

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere

✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because $\nabla_\theta \pi_\theta = 0$ on most of the state space. Hence**

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere
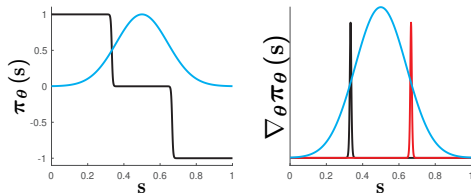
✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because $\nabla_\theta \pi_\theta = 0$ on most of the state space. Hence**

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere
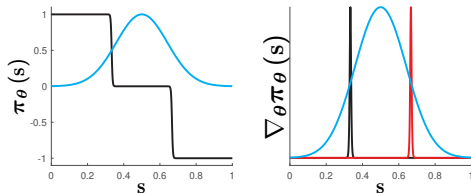
✓ Converge the policy to the true one over the learning

# RL & MPC for "strongly economic" problems

**Some policies are dominated by "switches", difficult to treat in RL because**
$\nabla_\theta \pi_\theta = 0$ **on most of the state space. Hence**

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

**is based on contributions from a very small number of samples. Parameter updates become "infrequent and jumpy".**

✓ Proposed policy relaxation techniques based on Interior-Point formulations, such that $\nabla_\theta \pi_\theta \neq 0$ almost everywhere

✓ Converge the policy to the true one over the learning

# Outline

# Reflections for today

**Focus on Economic problems**

- RLMPC is for performance
- Optimality "driven by external disturbances" seems the most interesting

# Reflections for today

**Focus on Economic problems**

- RLMPC is for performance
- Optimality "driven by external disturbances" seems the most interesting

**Not a competitor to other ideas**

- Keep classic approaches!
- Combinations are possible and beneficial
- RL-MPC "milks" the performance of other approaches

# Reflections for today

### RLMPC for constraint satisfaction

- Can "learn" to respect constraints
- Indirect approach, though
- ML-based "model-learning" better?

### Focus on Economic problems

- RLMPC is for performance
- Optimality "driven by external disturbances" seems the most interesting

### Not a competitor to other ideas

- Keep classic approaches!
- Combinations are possible and beneficial
- RL-MPC "milks" the performance of other approaches

# Reflections for today

**Focus on Economic problems**

- RLMPC is for performance
- Optimality "driven by external disturbances" seems the most interesting

**Not a competitor to other ideas**

- Keep classic approaches!
- Combinations are possible and beneficial
- RL-MPC "milks" the performance of other approaches

**RLMPC for constraint satisfaction**

- Can "learn" to respect constraints
- Indirect approach, though
- ML-based "model-learning" better?

**Software integration is a bottleneck**

- A lot of software for AI / RL
- Integration of MPC is not trivial

# Reflections for today

**RLMPC for constraint satisfaction**

- Can "learn" to respect constraints
- Indirect approach, though
- ML-based "model-learning" better?

**Focus on Economic problems**

- RLMPC is for performance
- Optimality "driven by external disturbances" seems the most interesting

**Software integration is a bottleneck**

- A lot of software for AI / RL
- Integration of MPC is not trivial

**Not a competitor to other ideas**

- Keep classic approaches!
- Combinations are possible and beneficial
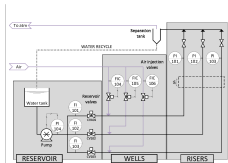- RL-MPC "milks" the performance of other approaches

**Prospects**:

- Software
- Stochastic constraints
- Dual mode / Optimized exploration
- Data efficiency
- Multi-agent problems, FATE
- More applications
- Can we make it a "technology"?

# Energy, Processes & Mobile robots

- Smart building
- Mobile robotics (UAV, USV)
- Wind energy
- Chemical process

- Smart house
- House with PV + Battery
- Energy Communities

Mix of experiments and simulations

**When does the best model fit produce the optimal policy?**
**I.e. when can we expect "classic MPC" to give us the highest performance?**

- Will do some repeats to put us in the right position to get there
- Introduce some "corollary" to the theory to explain our current understanding
- Show some basic examples

This is brand new lecture :-)

# Thanks for your attention!