Orientation

What we have seen:

• RLMPC, why does it work & different flavors

Orientation

What we have seen:

RLMPC, why does it work & different flavors

What we will do now:

- The theory is not about MPC, it is about MDPs (MPC is a special case)
- It has broader implications for AI and model-based decision making
- Provide (tentative) practical ideas for in Sim2Real

Al for Operational Decisions Some perspectives from recent research

Prof. Sebastien Gros

Department of Cybernetic
Faculty of Information Technology
NTNU

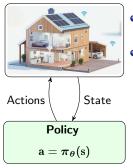
Freiburg PhD School

Outline

Decisions from data

2 Al models for Decision

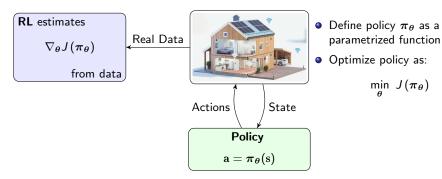




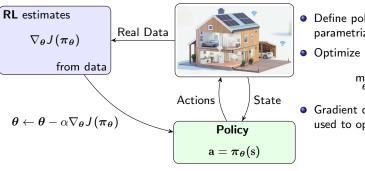
- Define policy π_{θ} as a parametrized function
- Optimize policy as:

$$\min_{\theta} J(\pi_{\theta})$$

4/18



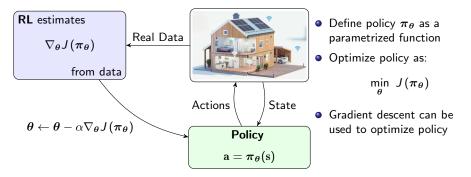
S. Gros (NTNU) Al for Decisions



- Define policy π_{θ} as a parametrized function
- Optimize policy as:

$$\min_{\theta} J(\pi_{\theta})$$

 Gradient descent can be used to optimize policy

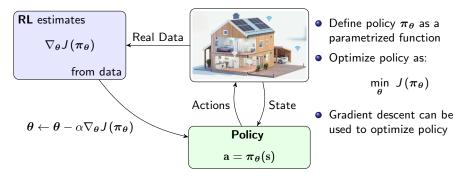


Remarks

- $\pi_{\theta}(s)$ often from DNN
- End-to-end training requires a lot of data and effort
- Difficult to provide explainability and guarantees on the policy we obtain
- Poor track record of industrial adoption (exception chatbots)

◆□▶ ◆□▶ ◆□▶ ◆□▶ ○□ ● ○○○○

4/18



Remarks

- $\pi_{\theta}(s)$ often from DNN
- End-to-end training requires a lot of data and effort
- Difficult to provide explainability and guarantees on the policy we obtain
- Poor track record of industrial adoption (exception chatbots)

→policy is often trained on high-fidelity simulations: embed knowledge, reduce costs, manage safety, promote explainability, control the training.

Real world: $\mathbf{s}_{k+1} \sim \rho(\cdot \mid \mathbf{s}_k, \mathbf{a}_k)$

Data with enough "richness"

fitting data

One-step model

Al model e.g.

$$\hat{\mathbf{s}}_{k+1} \sim \rho_{\boldsymbol{\theta}}(\cdot \,|\, \mathbf{s}_k, \mathbf{a}_k)$$

Real world: $s_{k+1} \sim \rho(\cdot | s_k, a_k)$

Data with enough "richness"

Here we no longer need the model ρ_{θ} to be "optimization-friendly"

S. Gros (NTNU)

Al for Decisions

Fall 2025

One-step model ${\bf Al\ model}\ {\bf e.g.}$ $\hat{\bf s}_{k+1} \sim \rho_{\boldsymbol \theta}(\cdot\,|\,{\bf s}_k,{\bf a}_k)$ fitting data

Real world: $\mathbf{s}_{k+1} \sim \rho(\cdot | \mathbf{s}_k, \mathbf{a}_k)$

Data with enough "richness"

5/18

Here we no longer need the model ρ_{θ} to be "optimization-friendly"

Model parameters θ

- Selected such that model ρ_{θ} "ressembles" reality ρ , using data (+ physics)
- Deterministic models are a special case of ρ_{θ}
- Many methods, from Least Squares and MLE to Bayesian and Adversarial Learning
- Still, in most cases ρ_{θ} "simplifies" ρ because
 - lacktriangle need for huge amounts of data to decide heta if model is very rich
 - lacktriangleright computationally demanding simulations if $ho_{m{ heta}}$ is expensive to sample from
- Then
 - biases in case of insufficiently rich model structure
 - validity is limited to some parts of the state-action space
 - ▶ only few first moments are correct (typ. mean + variance)

One-step model

Al model e.g.

$$\mathbf{\hat{s}}_{k+1} \sim
ho_{m{ heta}} ig(\cdot \, | \, \mathbf{s}_k, \mathbf{a}_k ig)$$
 fitting data

Virtual data from "simulating" $\hat{s}_{k+1} \sim \rho_{\theta}(\cdot\,|\,s_k,a_k) \text{ forward (Monte Carlo)}$

5/18

Here we no longer need the model ρ_{θ} to be "optimization-friendly"

Model parameters θ

- Selected such that model ρ_{θ} "ressembles" reality ρ , using data (+ physics)
- Deterministic models are a special case of $\rho_{m{ heta}}$
- Many methods, from Least Squares and MLE to Bayesian and Adversarial Learning
- Still, in most cases ρ_{θ} "simplifies" ρ because
 - lacktriangle need for huge amounts of data to decide $oldsymbol{ heta}$ if model is very rich
 - lacktriangleright computationally demanding simulations if $ho_{m{ heta}}$ is expensive to sample from
- Then
 - biases in case of insufficiently rich model structure
 - validity is limited to some parts of the state-action space
 - only few first moments are correct (typ. mean + variance)

One-step model

Al model e.g.

$$\mathbf{\hat{s}}_{k+1} \sim
ho_{m{ heta}} ig(\cdot \, | \, \mathbf{s}_k, \mathbf{a}_k ig)$$
 fitting data

Virtual data from "simulating" $\hat{s}_{k+1} \sim \rho_{\theta}(\cdot \, | \, s_k, a_k)$ forward (Monte Carlo)

In most cases, ρ_θ is a fairly simple representation of reality $\rho,$ easy to sample from

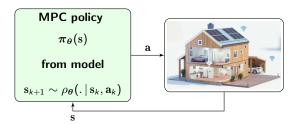
5/18

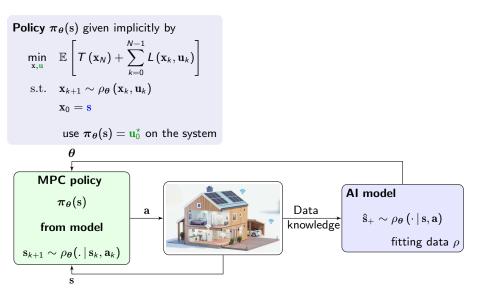
Here we no longer need the model ρ_{θ} to be "optimization-friendly"

Model parameters θ

- Selected such that model ρ_{θ} "ressembles" reality ρ , using data (+ physics)
- Deterministic models are a special case of ρ_{θ}
- Many methods, from Least Squares and MLE to Bayesian and Adversarial Learning
- Still, in most cases ρ_{θ} "simplifies" ρ because
 - lacktriangle need for huge amounts of data to decide $oldsymbol{ heta}$ if model is very rich
 - lacktriangledown computationally demanding simulations if $ho_{m{ heta}}$ is expensive to sample from
- Then
 - biases in case of insufficiently rich model structure
 - validity is limited to some parts of the state-action space
 - only few first moments are correct (typ. mean + variance)

$$\begin{split} \textbf{Policy} \ & \boldsymbol{\pi}_{\boldsymbol{\theta}}(\mathbf{s}) \ \text{given implicitly by} \\ & \underset{\mathbf{x},\mathbf{u}}{\text{min}} \quad \mathbb{E}\left[\left. \mathcal{T}\left(\mathbf{x}_{\mathit{N}}\right) + \sum_{k=0}^{\mathit{N}-1} L\left(\mathbf{x}_{\mathit{k}},\mathbf{u}_{\mathit{k}}\right) \right] \\ & \text{s.t.} \quad \mathbf{x}_{\mathit{k}+1} \sim \rho_{\boldsymbol{\theta}}\left(\mathbf{x}_{\mathit{k}},\mathbf{u}_{\mathit{k}}\right) \\ & \quad \mathbf{x}_{0} = \mathbf{s} \\ & \quad \text{use} \ \boldsymbol{\pi}_{\boldsymbol{\theta}}(\mathbf{s}) = \mathbf{u}_{0}^{\star} \ \text{on the system} \end{split}$$





Policy $\pi_{\theta}(s)$ given implicitly by

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

s.t.
$$\mathbf{x}_{k+1} \sim \rho_{\boldsymbol{\theta}} \left(\mathbf{x}_k, \mathbf{u}_k \right)$$

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

use $oldsymbol{\pi}_{oldsymbol{ heta}}(\mathbf{s}) = \mathbf{u}_0^{\star}$ on the system

Defines a paradigm...

- Performance from model accuracy
 - i.e. ρ_{θ} "close enough" to ρ
- "Ignore" that we replan all the time



- 4 ロ ト 4 個 ト 4 種 ト 4 種 ト - 種 - り Q ()

Policy $\pi_{\theta}(s)$ given implicitly by

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

s.t.
$$\mathbf{x}_{k+1} \sim \rho_{\boldsymbol{\theta}} \left(\mathbf{x}_k, \mathbf{u}_k \right)$$

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

use $\boldsymbol{\pi}_{\boldsymbol{\theta}}(\mathbf{s}) = \mathbf{u}_0^{\star}$ on the system

Relationship π_{θ} to π^* ??

- ullet They match if $ho_{m{ heta}}$ is exact and deterministic
- But it usually cannot be...



Policy $\pi_{\theta}(s)$ given implicitly by

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

s.t.
$$\mathbf{x}_{k+1} \sim \rho_{\boldsymbol{\theta}} \left(\mathbf{x}_k, \mathbf{u}_k \right)$$

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

use $oldsymbol{\pi}_{oldsymbol{ heta}}(\mathbf{s}) = \mathbf{u}_0^{\star}$ on the system

Relationship π_{θ} to π^{*} ??

- lacktriangle They match if $ho_{m{ heta}}$ is exact and deterministic
- But it usually cannot be...

"Standard methods" for choosing θ in general do not yield the best π_{θ} . RL over MPC can fix that, L becomes part of the model



Policy $\pi_{\theta}(\mathbf{s})$ given implicitly by $\min_{\mathbf{x},\mathbf{u}} \quad \mathbb{E}\left[T\left(\mathbf{x}_{N}\right) + \sum_{k=0}^{N-1} L\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)\right]$ s.t. $\mathbf{x}_{k+1} \sim \rho_{\theta}\left(\mathbf{x}_{k},\mathbf{u}_{k}\right)$ $\mathbf{x}_{0} = \mathbf{s}$

use $\pi_{\theta}(s) = \mathbf{u}_0^{\star}$ on the system

Relationship π_{θ} to π^* ??

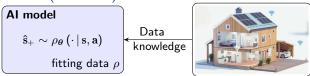
- \bullet They match if $\rho_{\pmb{\theta}}$ is exact and deterministic
- But it usually cannot be...

"Standard methods" for choosing θ in general do not yield the best π_{θ} . RL over MPC can fix that, L becomes part of the model



This is not the only way of taking decisions from models

Decision policy from Al models (Sim2Real)

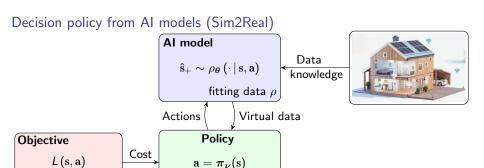


Classical Process:

• Fit Al model ρ_{θ} to real data

◆ロト ◆団 ト ◆ 豆 ト ◆ 豆 ・ り Q ()・

S. Gros (NTNU) Al for Decisions



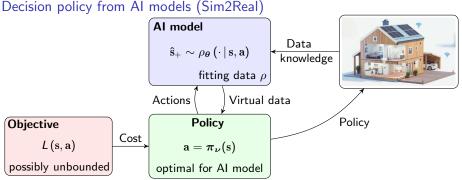
Classical Process:

- Fit Al model ρ_{θ} to real data
- ullet Develop optimal policy for L and $ho_{m{ heta}}$ from AI model, e.g. using In-Sim RL

optimal for AI model

- Define parametrized policy π_{ν}
- lacktriangle Optimize policy parameters $m{
 u}$ for performance w.r.t. Al model: $ho_{m{ heta}}, L
 ightarrow m{
 u}^{\star}$

S. Gros (NTNU) Al for Decisions

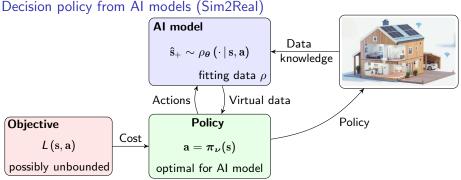


Classical Process:

- Fit Al model ρ_{θ} to real data
- ullet Develop optimal policy for L and ho_{ullet} from AI model, e.g. using In-Sim RL
 - Define parametrized policy $\pi_{
 u}$
 - lacktriangle Optimize policy parameters $m{
 u}$ for performance w.r.t. Al model: $ho_{m{ heta}}, L
 ightarrow m{
 u}^{\star}$
- ullet Transfer policy $\pi_{
 u^{\star}}$ into the real world

4□ > 4□ > 4□ > 4 = > □ 9 < ○</p>

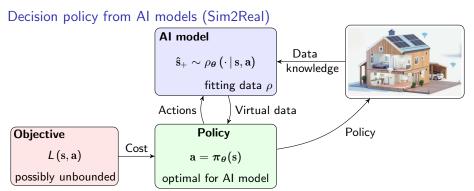
7/18



Classical Process:

- Fit AI model ρ_{θ} to real data
- Develop optimal policy for L and ρ_{θ} from AI model, e.g. using In-Sim RL
 - Define parametrized policy $\pi_{
 u}$
 - lacktriangle Optimize policy parameters $m{
 u}$ for performance w.r.t. Al model: $ho_{m{ heta}}, L
 ightarrow m{
 u}^{\star}$
- Transfer policy $\pi_{
 u^*}$ into the real world

Note: policy parameters u^* optimal for AI model become function of hetaLet's then label $\pi_{ heta}=\pi_{
u^*}$ for simplicity

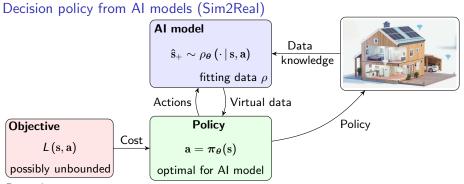


Remarks

- Relationship π_{θ} to π^* ??
 - ▶ They match if $\rho_{\theta} = \rho$, i.e. if model is exact
 - But ρ_{θ} is (almost always) an approximation

7/18

S. Gros (NTNU) Al for Decisions



Remarks

- Relationship π_{θ} to π^* ??
 - ▶ They match if $\rho_{\theta} = \rho$, i.e. if model is exact
 - ▶ But ρ_{θ} is (almost always) an approximation

Are "standard methods" for choosing θ resulting in a good policy π_{θ} ?? Empirically the answer is "no". Then how to choose θ to get a good policy?

4ロト 4個ト 4厘ト 4厘ト 厘 めQで

7/18

S. Gros (NTNU) Al for Decisions Fall 2025

The theory is about model-based decisions, and equivalences between MDPs

The theory is about model-based decisions, and equivalences between MDPs

World MDP

World MDP definition

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

S. Gros (NTNU)

The theory is about model-based decisions, and equivalences between MDPs

World MDP

World MDP definition

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

Optimal value functions

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

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

S. Gros (NTNU) Al for Decisions

The theory is about model-based decisions, and equivalences between MDPs

World MDP

World MDP definition

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

Optimal value functions

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

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

Optimal policy

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

S. Gros (NTNU) Al for Decisions

The theory is about model-based decisions, and equivalences between MDPs

World MDP

World MDP definition

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

Optimal value functions

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

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

Optimal policy

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

Model MDP (simulation environment)

Model MDP definition

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

Optimal value functions

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

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

Optimal policy

$$\hat{\pi}^{\star}\left(\mathbf{s}
ight) = \mathop{\mathsf{arg\,min}}\limits_{\mathbf{a}}\,\hat{Q}^{\star}\left(\mathbf{s},\mathbf{a}
ight)$$

The theory is about model-based decisions, and equivalences between MDPs

World MDP

World MDP definition

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

Optimal value functions

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

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

Optimal policy

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

Model MDP (simulation environment)

Model MDP definition

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

Optimal value functions

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

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

Optimal policy

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

4日ト 4日ト 4日ト 4日ト 日

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

The theory is about model-based decisions, and equivalences between MDPs

World MDP

World MDP definition

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

Optimal value functions

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

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

Optimal policy

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

Model MDP (simulation environment)

Model MDP definition

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

Optimal value functions

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

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

Optimal policy

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

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

<u>Proof</u>: telescopic sum, some non-trivial assumptions to prevent $\infty - \infty$ cancellations

More on MDP Equivalence

World MDP

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

Model MDP

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

S. Gros (NTNU) Al for Decisions

More on MDP Equivalence

World MDP

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

Model MDP

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

Theory says that

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

S. Gros (NTNU)

World MDP

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

Model MDP

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

Theory says that

- ullet Under some technical conditions there is a $\hat{oldsymbol{L}}$ such that $\hat{Q}^\star = Q^\star$
- ullet For $\hat{m{L}}=m{L}$ there is a (non-unique) "optimal" model $\hat{
 ho}$ such that $\hat{Q}^{\star}=Q^{\star}$

World MDP

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

Model MDP

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

Theory says that

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

World MDP

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

Model MDP

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

Theory says that

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

S. Gros (NTNU)

World MDP

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

Model MDP

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

Theory says that

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

In-Sim policy training uses a Model MDP

- ullet Practitioners mostly work on model $\hat{
 ho}$, but "unsure" on how to tune it
- In the ML communities, people talk about "value alignment", we are trying to make this "MDP equivalence" understood

World MDP

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

Model MDP

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

Theory says that

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

Remarks

• Non-unique optimal model $\hat{\rho}$ leaves room for aligning it to the real world (classical fitting)



World MDP

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

Model MDP

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

Theory says that

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

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it to the real world (classical fitting)
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then we can have support $\hat{\rho} \subset$ support of ρ

S. Gros (NTNU)

World MDP

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

Model MDP

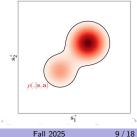
- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that

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

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it to the real world (classical fitting)
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then we can have support $\hat{\rho} \subset$ support of ρ



World MDP

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

Model MDP

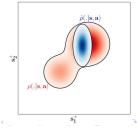
- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that

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

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it to the real world (classical fitting)
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then we can have support $\hat{\rho} \subset$ support of ρ



9/18

World MDP

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

Model MDP

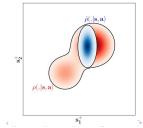
- States s and actions a
- Cost $\hat{L}(s, a)$
- ullet Transition $\mathbf{s}_+ \sim \hat{oldsymbol{
 ho}} \left[\cdot | \mathbf{s}, \mathbf{a}
 ight]$

Theory says that

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

Remarks

- Non-unique optimal model $\hat{\rho}$ leaves room for aligning it to the real world (classical fitting)
- If V^* is continuous and support of ρ is bounded(?) and path-connected, then we can have support $\hat{\rho} \subset$ support of ρ
- More simply said: we can build optimal models $\hat{\rho}$ that make "plausible" predictions about the real world.

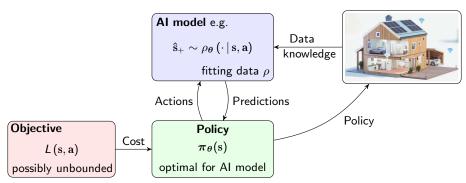


Outline

Decisions from data

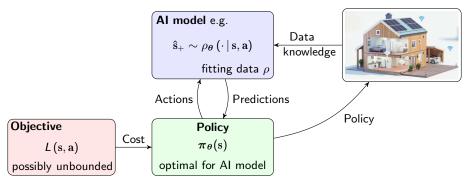
2 Al models for Decision





Learning model ρ_{θ} should aim at identifying

$$\theta^{\star} = \underset{\theta}{\operatorname{arg\,min}} \ J(\pi_{\theta}) \tag{1}$$

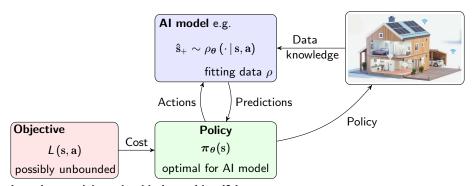


Learning model ρ_{θ} should aim at identifying

$$\theta^{\star} = \underset{\theta}{\operatorname{arg\,min}} \ J\left(\pi_{\theta}\right) \tag{1}$$

ullet Closed-loop performance J is intricate (heta o simulations o policy o real system)

◆ロト ◆個ト ◆意ト ◆意ト ・ 意 ・ かへで

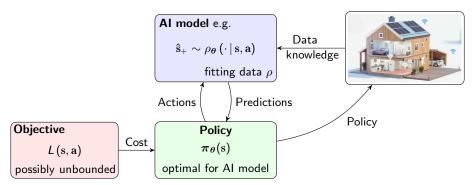


Learning model ρ_{θ} should aim at identifying

$$\theta^* = \underset{\theta}{\operatorname{arg\,min}} \ J(\pi_{\theta}) \tag{1}$$

- ullet Closed-loop performance J is intricate $(oldsymbol{ heta} o$ simulations o policy o real system)
- ullet $\nabla_{ heta} J(\pi_{ heta})$ is to be estimated from data, i.e. data hungry, noisy, possible biases

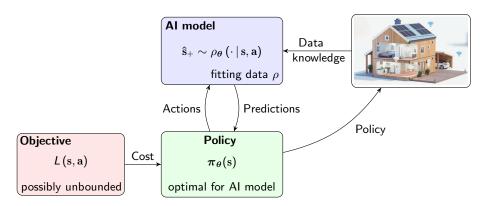
◆ロト ◆部 ト ◆ 差 ト ◆ 差 ・ 夕 Q (*)



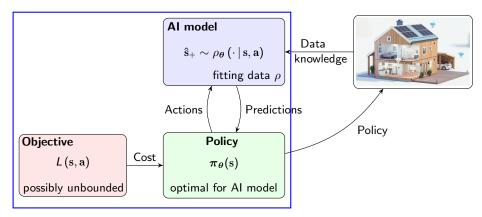
Learning model ρ_{θ} should aim at identifying

$$\theta^* = \underset{\theta}{\operatorname{arg\,min}} \ J(\pi_{\theta}) \tag{1}$$

- ullet Closed-loop performance J is intricate $(oldsymbol{ heta} o$ simulations o policy o real system)
- $\nabla_{\theta} J(\pi_{\theta})$ is to be estimated from data, i.e. data hungry, noisy, possible biases
- (1) is in general different than model fitting, i.e. no loss function does (1)

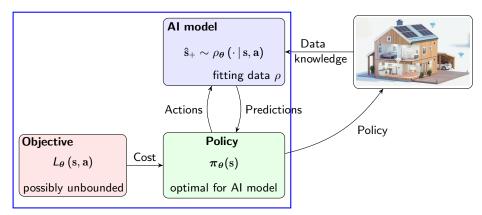


◆ロト ◆団 ト ◆ 豆 ト ◆ 豆 ・ り Q ()・



• The model is not ρ_{θ} , it is the entire "decision-making box"

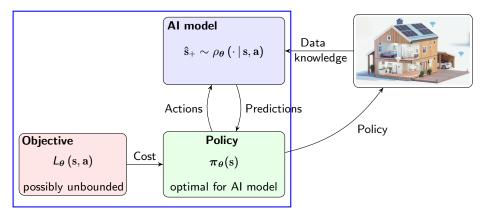
◆ロト ◆個 ト ◆ 恵 ト ◆ 恵 ・ からで



- The model is not ρ_{θ} , it is the entire "decision-making box"
- The model includes objective L_{θ} used to build the policy*

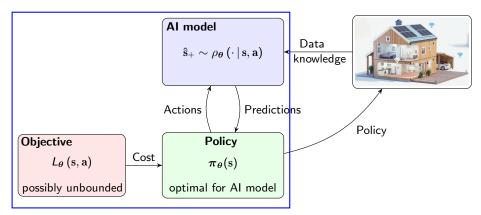
◆ロト ◆問 ト ◆ 恵 ト ◆ 恵 ・ 釣 ♀ ○○

^{*} policy performance on real system still assessed via L



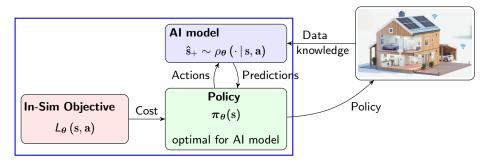
- The model is not ρ_{θ} , it is the entire "decision-making box"
- The model includes objective L_{θ} used to build the policy*
- Best model ρ_{θ} should not necessarily represent the data in a "classical sense"

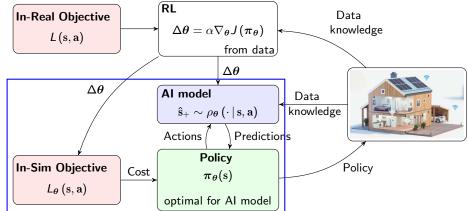
^{*} policy performance on real system still assessed via L

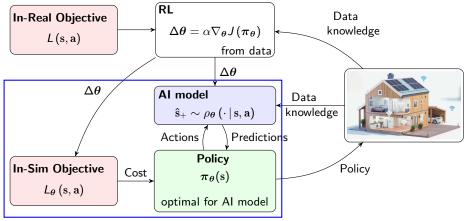


Theorem: under some (technical) assumptions, there is a L_{θ} such that $\pi_{\theta}=\pi_{\star}$, even if ρ_{θ} does not represent real world ρ correctly

◆ロ → ◆母 → ◆ き → ◆ き → り へ ○





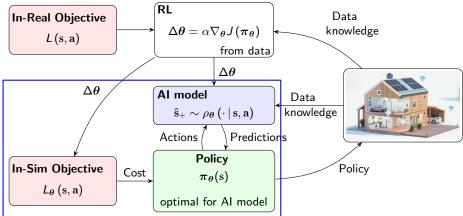


Policy gradient

$$abla_{ heta} J(\pi_{ heta}) = \mathbb{E}\left[
abla_{ heta} \pi_{ heta}
abla_{ ext{a}} Q^{\pi_{ heta}}
ight]$$

- $Q^{\pi_{\theta}}$ is the critic, well-established RL tool
- $\nabla_{\theta} \pi_{\theta}$ requires differentiating the closed-loop simulation...

◆□▶◆□▶◆□▶◆□▶ □ めの@



Difficulty: computing $\nabla_{\theta} \pi_{\theta}$ requires total differentiation through the simulation-based policy optimization process:

 $oldsymbol{ heta}
ightarrow \operatorname{simulations}
ightarrow \operatorname{optimal policy}$ for $ho_{oldsymbol{ heta}}$

Differentiating through policy optimization can be computationally heavy

Performance for a stochastic policy π_{ν} , model MDP defined by L_{θ} , $\hat{\rho}_{\theta}$, γ

$$\hat{J}_{\boldsymbol{\theta}}\left(\pi_{\boldsymbol{\nu}}\right) = \mathbb{E}_{\hat{\boldsymbol{\theta}}_{\boldsymbol{\theta}}}\left[\left.\sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right)\right| \mathbf{a}_{k} \sim \pi_{\boldsymbol{\nu}}\left[\cdot | \mathbf{s}_{k}\right]\right]$$

Performance for a stochastic policy π_{ν} , model MDP defined by L_{θ} , $\hat{\rho}_{\theta}$, γ

$$\hat{J}_{\boldsymbol{\theta}}\left(\pi_{\boldsymbol{\nu}}\right) = \mathbb{E}_{\hat{\boldsymbol{\theta}}_{\boldsymbol{\theta}}}\left[\left.\sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right)\right| \mathbf{a}_{k} \sim \pi_{\boldsymbol{\nu}}\left[\cdot | \mathbf{s}_{k}\right]\right]$$

Policy given by

$$\mathbf{\nu}^{\star} = \underset{\mathbf{\nu}}{\operatorname{arg\,min}} \ \hat{J}_{\boldsymbol{\theta}}\left(\pi_{\nu}\right) \qquad \text{or equivalently} \qquad \frac{\partial}{\partial \nu} \hat{J}_{\boldsymbol{\theta}}\left(\pi_{\nu}\right) = 0$$

◆ロト ◆部 ト ◆ 恵 ト ◆ 恵 ・ 釣 へ ご

Performance for a stochastic policy π_{ν} , model MDP defined by L_{θ} , $\hat{\rho}_{\theta}$, γ

$$\hat{J}_{\boldsymbol{\theta}}\left(\pi_{\boldsymbol{\nu}}\right) = \mathbb{E}_{\hat{\boldsymbol{\theta}}_{\boldsymbol{\theta}}}\left[\left.\sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right)\right| \mathbf{a}_{k} \sim \pi_{\boldsymbol{\nu}}\left[\cdot | \mathbf{s}_{k}\right]\right]$$

Policy given by

$$\mathbf{\nu}^{\star} = \underset{\mathbf{\nu}}{\operatorname{arg \, min}} \ \hat{J}_{\boldsymbol{\theta}} \left(\pi_{\boldsymbol{\nu}} \right) \qquad \text{or equivalently} \qquad \frac{\partial}{\partial \boldsymbol{\nu}} \hat{J}_{\boldsymbol{\theta}} \left(\pi_{\boldsymbol{\nu}} \right) = 0$$

Policy sensitivity: if we change θ (cost L_{θ} and simulation $\hat{\rho}_{\theta}$) how does π_{ν} change?

$$\frac{\partial^{2}}{\partial \nu^{2}}\hat{J}_{\theta}\left(\pi_{\nu}\right)\underbrace{\frac{\mathrm{d}\pi_{\nu}}{\mathrm{d}\theta}}_{=\nabla_{\theta}\pi_{\theta}}+\frac{\partial^{2}}{\partial \nu\partial\theta}\hat{J}_{\theta}\left(\pi_{\nu}\right)=0$$

Performance for a stochastic policy π_{ν} , model MDP defined by L_{θ} , $\hat{\rho}_{\theta}$, γ

$$\hat{J}_{\boldsymbol{\theta}}\left(\pi_{\boldsymbol{\nu}}\right) = \mathbb{E}_{\hat{\boldsymbol{\theta}}_{\boldsymbol{\theta}}}\left[\left.\sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right)\right| \, \mathbf{a}_{k} \sim \pi_{\boldsymbol{\nu}}\left[\cdot | \mathbf{s}_{k}\right]\right]$$

Policy given by

$$\mathbf{v}^{\star} = \underset{\mathbf{v}}{\operatorname{arg min}} \ \hat{J}_{\theta}\left(\pi_{\nu}\right) \qquad \text{or equivalently} \qquad \frac{\partial}{\partial \nu} \hat{J}_{\theta}\left(\pi_{\nu}\right) = 0$$

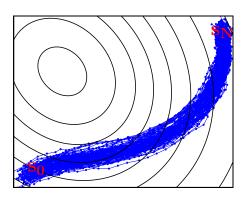
Policy sensitivity: if we change θ (cost L_{θ} and simulation $\hat{\rho}_{\theta}$) how does π_{ν} change?

$$\frac{\partial^{2}}{\partial \nu^{2}} \hat{J}_{\theta} (\pi_{\nu}) \underbrace{\frac{\mathrm{d}\pi_{\nu}}{\mathrm{d}\theta}}_{=\nabla_{\theta}\pi_{\theta}} + \frac{\partial^{2}}{\partial \nu \partial \theta} \hat{J}_{\theta} (\pi_{\nu}) = 0$$

Simulation: (with $\mathbf{s} \equiv \mathbf{s}_{0,\dots,N}$ and $\mathbf{a} \equiv \mathbf{a}_{0,\dots,N}$)

$$\hat{J}_{\boldsymbol{\theta}}\left(\pi_{\nu}\right) = \int \sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right) \varphi\left(\mathbf{s}, \mathbf{a}\right) d\mathbf{s} d\mathbf{a}$$

$$\varphi\left(\mathbf{s},\mathbf{a}\right) = \rho_0\left[\mathbf{s}_0\right] \prod_{k=0}^{N-1} \hat{\rho}_{\boldsymbol{\theta}}\left[\mathbf{s}_{k+1} | \mathbf{s}_k, \mathbf{a}_k\right] \prod_{k=0}^{N} \pi_{\boldsymbol{\nu}}\left[\mathbf{a}_k | \mathbf{s}_k\right] \qquad \text{(density of the simulation)}$$



• Blue trajectories are realization of φ :

$$\mathbf{s}^{i},\mathbf{a}^{i}\simarphi\left(\cdot,\cdot
ight)$$

Sample-based evaluation (n samples)

$$\hat{J}_{\boldsymbol{\theta}}\left(\pi_{\nu}\right) pprox rac{1}{n} \sum_{k=i}^{n} \sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}^{i}, \mathbf{a}_{k}^{i}\right)$$

Simulation: (with $\mathbf{s} \equiv \mathbf{s}_{0,\dots,N}$ and $\mathbf{a} \equiv \mathbf{a}_{0,\dots,N}$)

$$\hat{J}_{\theta}(\pi_{\nu}) = \int \sum_{k=0}^{N} \gamma^{k} L_{\theta}(\mathbf{s}_{k}, \mathbf{a}_{k}) \varphi(\mathbf{s}, \mathbf{a}) \, d\mathbf{s} d\mathbf{a}$$

$$\varphi(\mathbf{s}, \mathbf{a}) = \rho_0[\mathbf{s}_0] \prod_{k=0}^{N-1} \hat{\rho}_{\boldsymbol{\theta}}[\mathbf{s}_{k+1} | \mathbf{s}_k, \mathbf{a}_k] \prod_{k=0}^{N} \pi_{\boldsymbol{\nu}}[\mathbf{a}_k | \mathbf{s}_k]$$

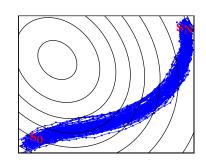
(density of the simulation)

◆ロト ◆昼 ト ◆ 恵 ト ・ 恵 ・ かへで

Simulation

$$\hat{J}_{\theta}(\pi_{\nu}) = \int \sum_{k=0}^{N} \gamma^{k} L_{\theta}(\mathbf{s}_{k}, \mathbf{a}_{k}) \varphi(\mathbf{s}, \mathbf{a}) \, d\mathbf{s} d\mathbf{a}$$

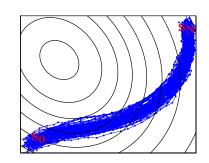
$$\varphi(\mathbf{s}, \mathbf{a}) = \rho_{0} \left[\mathbf{s}_{0}\right] \prod_{k=0}^{N-1} \hat{\rho}_{\theta} \left[\mathbf{s}_{k+1} | \mathbf{s}_{k}, \mathbf{a}_{k}\right] \prod_{k=0}^{N} \pi_{\nu} \left[\mathbf{a}_{k} | \mathbf{s}_{k}\right]$$



15 / 18

Simulation

$$\begin{split} \hat{J}_{\boldsymbol{\theta}}\left(\boldsymbol{\pi}_{\boldsymbol{\nu}}\right) &= \int \sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right) \boldsymbol{\varphi}\left(\mathbf{s}, \mathbf{a}\right) \mathrm{d} \mathbf{s} \mathrm{d} \mathbf{a} \\ \boldsymbol{\varphi}\left(\mathbf{s}, \mathbf{a}\right) &= \rho_{0} \left[\mathbf{s}_{0}\right] \prod_{k=0}^{N-1} \hat{\rho}_{\boldsymbol{\theta}} \left[\mathbf{s}_{k+1} | \mathbf{s}_{k}, \mathbf{a}_{k}\right] \prod_{k=0}^{N} \pi_{\boldsymbol{\nu}} \left[\mathbf{a}_{k} | \mathbf{s}_{k}\right] \end{split}$$



Direct Differentiation:

$$\partial \hat{J}_{\boldsymbol{\theta}}\left(\pi_{\boldsymbol{\nu}}\right) = \mathbb{E}_{\hat{\boldsymbol{\theta}}_{\boldsymbol{\theta}}}\left[\left.\sum_{k=0}^{N} \gamma^{k} \left(\partial L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right) + L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right) \partial \log \varphi\left(\mathbf{s}, \mathbf{a}\right)\right)\right| \, \mathbf{a}_{k} \sim \pi_{\boldsymbol{\nu}}\left[\cdot | \mathbf{s}_{k} \right]\right]$$

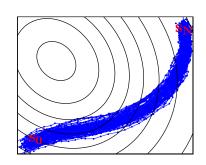
can be evaluated from data

◆ロト ◆母 ト ◆ 差 ト ◆ 差 ・ 釣 へ ②

15 / 18

Simulation

$$\begin{split} \hat{J}_{\boldsymbol{\theta}}\left(\pi_{\boldsymbol{\nu}}\right) &= \int \sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right) \varphi\left(\mathbf{s}, \mathbf{a}\right) \mathrm{d}\mathbf{s} \mathrm{d}\mathbf{a} \\ \varphi\left(\mathbf{s}, \mathbf{a}\right) &= \rho_{0} \left[\mathbf{s}_{0}\right] \prod_{k=0}^{N-1} \hat{\rho}_{\boldsymbol{\theta}} \left[\mathbf{s}_{k+1} | \mathbf{s}_{k}, \mathbf{a}_{k}\right] \prod_{k=0}^{N} \pi_{\boldsymbol{\nu}} \left[\mathbf{a}_{k} | \mathbf{s}_{k}\right] \end{split}$$



Direct Differentiation:

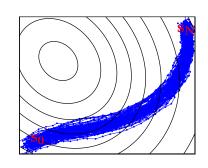
$$\partial \hat{J}_{\theta} \left(\pi_{\nu} \right) \approx \frac{1}{n} \sum_{i=1}^{n} \sum_{k=0}^{N} \gamma^{k} \left(\partial L_{\theta} \left(\mathbf{s}_{k}^{i}, \mathbf{a}_{k}^{i} \right) + L_{\theta} \left(\mathbf{s}_{k}^{i}, \mathbf{a}_{k}^{i} \right) \partial \log \varphi \left(\mathbf{s}_{k}^{i}, \mathbf{a}_{k}^{i} \right) \right)$$

◆ロト ◆母 ト ◆ 差 ト ◆ 差 ・ 釣 へ ②

15 / 18

Simulation

$$\begin{split} \hat{J}_{\boldsymbol{\theta}}\left(\pi_{\boldsymbol{\nu}}\right) &= \int \sum_{k=0}^{N} \gamma^{k} L_{\boldsymbol{\theta}}\left(\mathbf{s}_{k}, \mathbf{a}_{k}\right) \varphi\left(\mathbf{s}, \mathbf{a}\right) \mathrm{d}\mathbf{s} \mathrm{d}\mathbf{a} \\ \varphi\left(\mathbf{s}, \mathbf{a}\right) &= \rho_{0} \left[\mathbf{s}_{0}\right] \prod_{k=0}^{N-1} \hat{\rho}_{\boldsymbol{\theta}} \left[\mathbf{s}_{k+1} | \mathbf{s}_{k}, \mathbf{a}_{k}\right] \prod_{k=0}^{N} \pi_{\boldsymbol{\nu}} \left[\mathbf{a}_{k} | \mathbf{s}_{k}\right] \end{split}$$



Direct Differentiation:

$$\partial \hat{J}_{\boldsymbol{\theta}} \left(\pi_{\boldsymbol{\nu}} \right) \approx \frac{1}{n} \sum_{i=1}^{n} \sum_{k=0}^{N} \gamma^{k} \left(\partial L_{\boldsymbol{\theta}} \left(\mathbf{s}_{k}^{i}, \mathbf{a}_{k}^{i} \right) + L_{\boldsymbol{\theta}} \left(\mathbf{s}_{k}^{i}, \mathbf{a}_{k}^{i} \right) \partial \log \varphi \left(\mathbf{s}_{k}^{i}, \mathbf{a}_{k}^{i} \right) \right)$$

Remarks

- √ Second-order sensitivities are similar
- √ Simple and computationally very efficient
- √ Can differentiate through discrete state and action sets
- Sample-based estimations are very noisy

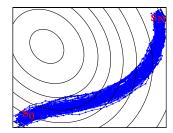


15 / 18

In-Sim Actor-critic does

$$\nabla_{\nu}\hat{J}_{\pmb{\theta}}\left(\pi_{\nu}\right) = \mathbb{E}_{\begin{array}{c} s \sim \rho_{\pmb{\theta}} \\ a \sim \pi_{\nu} \end{array}} \left[\nabla_{\nu}\log\pi_{\nu}\left(a|s\right)A^{\pi_{\nu}}\left(s,a\right)\right] = 0$$

defines relationship $\pi_{\boldsymbol{\nu}^{\star}}$ with $\boldsymbol{\nu}^{\star}$ function of $\boldsymbol{\theta}$



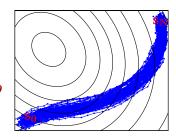
In-Sim Actor-critic does

$$\nabla_{\nu}\hat{J}_{\boldsymbol{\theta}}\left(\pi_{\nu}\right) = \mathbb{E}_{\substack{\mathbf{s} \sim \rho_{\boldsymbol{\theta}} \\ \mathbf{a} \sim \pi_{\nu}}} \left[\nabla_{\nu}\log \pi_{\nu}\left(\mathbf{a}|\mathbf{s}\right)\mathcal{A}^{\pi_{\nu}}\left(\mathbf{s},\mathbf{a}\right)\right] = 0$$

defines relationship π_{ν^*} with ν^* function of θ

In-Real Actor-critic evaluates

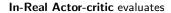
$$\nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\nu}^{\star}}) = \mathbb{E}_{\mathsf{Real World}} \left[\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\nu}^{\star}} \left(\mathbf{a} | \mathbf{s} \right) A^{\boldsymbol{\nu}^{\star}} \left(\mathbf{s}, \mathbf{a} \right) \right]$$



In-Sim Actor-critic does

$$abla_{
u}\hat{J_{m{ heta}}}\left(\pi_{
u}
ight) = \mathbb{E} egin{array}{l} \mathbf{s} \sim
ho_{m{ heta}} \ \mathbf{a} \sim \pi_{
u} \end{array} \left[
abla_{
u} \log \pi_{
u} \left(\mathbf{a} | \mathbf{s}
ight) \mathcal{A}^{m{\pi}_{
u}} \left(\mathbf{s}, \mathbf{a}
ight) = 0$$

defines relationship $\pi_{\boldsymbol{\nu}^{\star}}$ with $\boldsymbol{\nu}^{\star}$ function of $\boldsymbol{\theta}$

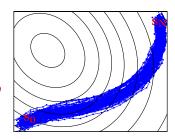


$$abla_{m{ heta}} J\left(\pi_{m{
u}^{\star}}
ight) = \mathbb{E}_{\mathsf{Real\ World}}\left[
abla_{m{ heta}} \log \pi_{m{
u}^{\star}}\left(\mathbf{a}|\mathbf{s}
ight) m{\mathcal{A}}^{m{
u}^{\star}}\left(\mathbf{s},\mathbf{a}
ight)
ight]$$

where $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\nu}^{\star}} \left(\mathbf{a} | \mathbf{s} \right)$ requires $\frac{\mathrm{d} \boldsymbol{\nu}^{\star}}{\mathrm{d} \boldsymbol{\theta}}$



$$\frac{\partial^{2}}{\partial \boldsymbol{\nu}^{2}} \hat{J}_{\boldsymbol{\theta}} \left(\pi_{\boldsymbol{\nu}^{\star}} \right) \frac{\mathrm{d} \boldsymbol{\nu}^{\star}}{\mathrm{d} \boldsymbol{\theta}} + \frac{\partial^{2}}{\partial \boldsymbol{\nu} \partial \boldsymbol{\theta}} \hat{J}_{\boldsymbol{\theta}} \left(\pi_{\boldsymbol{\nu}^{\star}} \right) = 0$$

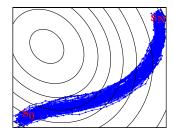


16 / 18

In-Sim Actor-critic does

$$abla_{
u}\hat{J_{m{ heta}}}\left(\pi_{
u}
ight) = \mathbb{E} egin{array}{c} \mathbf{s} \sim
ho_{m{ heta}} \ \mathbf{a} \sim \pi_{
u} \end{array} \left[
abla_{
u} \log \pi_{
u} \left(\mathbf{a} | \mathbf{s}
ight) \mathcal{A}^{m{\pi}_{
u}} \left(\mathbf{s}, \mathbf{a}
ight) = 0$$

defines relationship $\pi_{\boldsymbol{\nu}^*}$ with $\boldsymbol{\nu}^*$ function of $\boldsymbol{\theta}$



In-Real Actor-critic evaluates

$$\nabla_{\boldsymbol{\theta}} J\left(\pi_{\boldsymbol{\nu}^{\star}}\right) = \mathbb{E}_{\mathsf{Real World}}\left[\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\nu}^{\star}}\left(\mathbf{a} | \mathbf{s}\right) \boldsymbol{\mathcal{A}}^{\boldsymbol{\nu}^{\star}}\left(\mathbf{s}, \mathbf{a}\right)\right]$$

where $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\nu}^{\star}} \left(\mathbf{a} | \mathbf{s} \right)$ requires $\frac{\mathrm{d} \boldsymbol{\nu}^{\star}}{\mathrm{d} \boldsymbol{\theta}}$

Sensitivity

$$\frac{\partial^{2}}{\partial \boldsymbol{\nu}^{2}}\hat{J}_{\boldsymbol{\theta}}\left(\boldsymbol{\pi}_{\boldsymbol{\nu}^{\star}}\right)\frac{\mathrm{d}\boldsymbol{\nu}^{\star}}{\mathrm{d}\boldsymbol{\theta}}+\frac{\partial^{2}}{\partial\boldsymbol{\nu}\partial\boldsymbol{\theta}}\hat{J}_{\boldsymbol{\theta}}\left(\boldsymbol{\pi}_{\boldsymbol{\nu}^{\star}}\right)=0$$

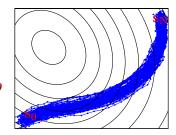
Idea: compute $\frac{\mathrm{d}\nu^{\star}}{\mathrm{d}\theta}$ from differentiating the in-sim actor-critic w.r.t. $\frac{\partial}{\partial\nu}$ and $\frac{\partial}{\partial\theta}$

16 / 18

In-Sim Actor-critic does

$$abla_{
u}\hat{J_{m{ heta}}}\left(\pi_{
u}
ight) = \mathbb{E} egin{array}{c} \mathbf{s} \sim
ho_{m{ heta}} \ \mathbf{a} \sim \pi_{
u} \end{array} \left[
abla_{
u} \log \pi_{
u} \left(\mathbf{a} | \mathbf{s}
ight) \mathcal{A}^{m{\pi}_{
u}} \left(\mathbf{s}, \mathbf{a}
ight) = 0$$

defines relationship $\pi_{\boldsymbol{\nu}^*}$ with $\boldsymbol{\nu}^*$ function of $\boldsymbol{\theta}$



In-Real Actor-critic evaluates

$$abla_{m{ heta}} J\left(\pi_{m{
u}^{\star}}
ight) = \mathbb{E}_{\mathsf{Real World}}\left[
abla_{m{ heta}} \log \pi_{m{
u}^{\star}}\left(\mathbf{a}|\mathbf{s}
ight) m{A}^{m{
u}^{\star}}\left(\mathbf{s},\mathbf{a}
ight)
ight]$$

where $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\nu}^*} \left(\mathbf{a} | \mathbf{s} \right)$ requires $\frac{\mathrm{d} \boldsymbol{\nu}^*}{\mathrm{d} \boldsymbol{\theta}}$

Sensitivity

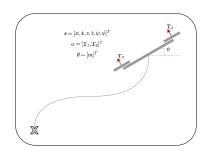
$$\frac{\partial^{2}}{\partial \boldsymbol{\nu}^{2}}\hat{J}_{\boldsymbol{\theta}}\left(\boldsymbol{\pi}_{\boldsymbol{\nu}^{\star}}\right)\frac{\mathrm{d}\boldsymbol{\nu}^{\star}}{\mathrm{d}\boldsymbol{\theta}}+\frac{\partial^{2}}{\partial \boldsymbol{\nu}\partial \boldsymbol{\theta}}\hat{J}_{\boldsymbol{\theta}}\left(\boldsymbol{\pi}_{\boldsymbol{\nu}^{\star}}\right)=0$$

Idea: compute $\frac{d\nu^*}{d\theta}$ from differentiating the in-sim actor-critic w.r.t. $\frac{\partial}{\partial \nu}$ and $\frac{\partial}{\partial \theta}$

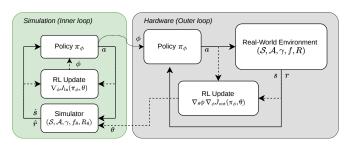
Proposed first algorithm to do that.

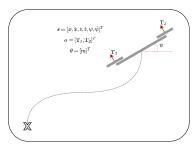
Computationally heavy, but we have "missed" some simplifications, testing further

Example - In-Real RL over In-Sim RL

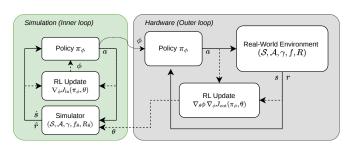


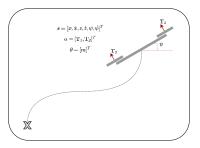
Example - In-Real RL over In-Sim RL

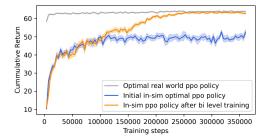




Example - In-Real RL over In-Sim RL







◆□▶ ◆圖▶ ◆意▶ ◆意▶

Norwegian Center on AI for Decisions



Email



- Start-up phase
- Recruiting soon!
- We can apply for funding to bring strong international postdocs



Thank you!