



Safe model-based multi-agent mean-field reinforcement learning

Conference Poster

Author(s):

Jusup, Matej ; Pásztor, Barna; Janik, Tadeusz; Zhang, Kenan; Corman, Francesco ; Krause, Andreas; Bogunovic, Ilija

Publication date:

2024-05

Permanent link:

<https://doi.org/10.3929/ethz-b-000673391>

Rights / license:

In Copyright - Non-Commercial Use Permitted

Funding acknowledgement:

181210 - DADA - Dynamic data driven Approaches for stochastic Delay propagation Avoidance in railways (SNF)

Safe Model-Based Multi-Agent Mean-Field Reinforcement Learning

M. Jusup¹, B. Pasztor¹, T. Janik¹, K. Zhang², F. Corman¹, A. Krause¹, I. Bogunovic³
¹ETH Zurich; ²EPFL Lausanne; ³UCL London



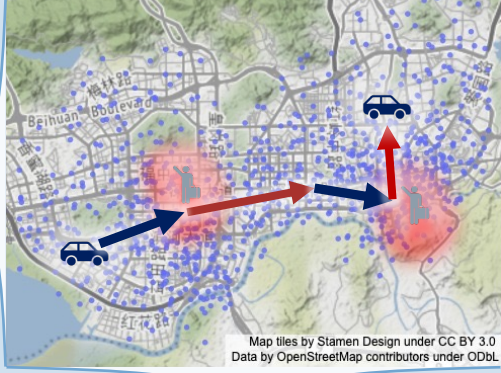
Motivation: Vehicle Repositioning

Repositioning to the high-demand areas

\$\$\$

- Fair service accessibility
- Limit on number of vehicles per district

Safety constraints!



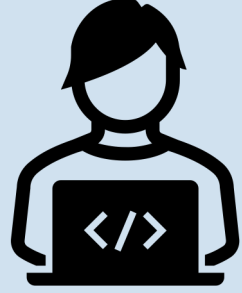
- True mean-field distribution: $\mu_{n,t}$
- True safety constraints: $h_C(\mu_{n,t}) \geq 0$

Safe Model-Based Mean-Field RL

RA collects trajectories $(z_{n,t}, s_{n,t+1})$

$h_C(\tilde{\mu}_{n,t}) \geq 0 \Rightarrow h_C(\mu_{n,t}) \geq 0?$

Controller sends updated policy π_n^* to RA



- Statistical model: $\tilde{f}_{n-1}(z) \approx m_{n-1}(z) + \Sigma_{n-1}(z)$
- Statistical mean-field distribution: $\tilde{\mu}_{n,t}$
- Statistical safety constraints: $h_C(\tilde{\mu}_{n,t}) \geq 0$
- Policy: π_n

Why Safe Mean-Field Reinforcement Learning?

Objective



Learning the optimal *safe* policy π^* under *unknown* transitions f given *safety* constraints $h_C(\mu_{n,t}) \geq 0$

Individual interactions lead to the combinatorial state-action space

Mean-field distribution μ of cooperative identical agents

The representative agent (RA) interacts with the mean-field distribution μ

RA policy π^* is used to control all the agents

-  No individual interactions
-  Complex inputs $z = (s, \mu, a)$ and probabilistic transitions $U(\cdot)$

Safe-M³-UCRL

$$\pi_n^* = \arg \max_{\pi_n \in \Pi} \max_{\eta(\cdot) \in [-1, 1]^p} \mathbb{E} \left[\sum_{t=0}^{T-1} r(\tilde{z}_{n,t}) \mid \tilde{\mu}_{n,0} = \mu_0 \right]$$

subject to $\tilde{a}_{n,t} = \pi_{n,t}(\tilde{s}_{n,t}, \tilde{\mu}_{n,t})$

$$\tilde{f}_{n-1}(\tilde{z}_{n,t}) = \mathbf{m}_{n-1}(\tilde{z}_{n,t}) + \beta_{n-1} \Sigma_{n-1}(\tilde{z}_{n,t}) \eta(\tilde{z}_{n,t})$$

$$\tilde{s}_{n,t+1} = \tilde{f}_{n-1}(\tilde{z}_{n,t}) + \varepsilon_{n,t}$$

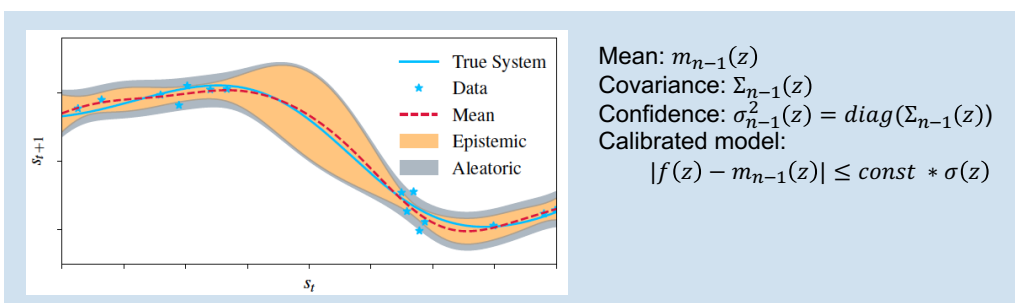
$$\tilde{\mu}_{n,t+1} = U(\tilde{\mu}_{n,t}, \pi_{n,t}, \tilde{f}_{n-1})$$

$h_C(\tilde{\mu}_{n,t+1}) \geq L_h C_{n,t+1}$ ← Enables safe exploration!

Contributions of Safe-M³-UCRL

- Safe exploration guided by budget $C_{n,t}$ induced by epistemic uncertainty $\sigma_{n-1}(z)$
- Relationship between safety constraints under statistical and true environments $|h_C(\tilde{\mu}_{n,t}) - h_C(\mu_{n,t})| \leq L_h C_{n,t}$
- Algorithm that adheres to the safety constraints throughout the entire execution (with high probability)
- Showcasing usefulness of Mean-Field RL in real-world applications!**

Calibrated Statistical Model of Unknown Transitions



Model-Based Learning Protocol in Safe-M³-UCRL

Input: Safety constraint $h_C(\cdot)$, initial mean-field distribution μ_0 , number of episodes N , number of steps T

- for $n = 1, \dots, N$ do
- Compute $C_{n,t}$ for $t = 1, \dots, T$
- Learn a policy π_n^* by optimizing the objective
- Execute the obtained policy π_n^*
- Collect the trajectories from the representative agent
- Update the statistical model \tilde{f}_{n-1}
- end for

Return π_N^*

