

AIMLC ZG512 -
Deep Reinforcement Learning


 **BITS Pilani**
Pursuing Global Education

Session 17: Special Topics in DRL

Introducing - Safety in Reinforcement Learning


S. P. Vimal, Department of CSIS, WILP Division (vimalsp@wilp.bits-pilani.ac.in)

1

 **Safety**

BIG QUESTION:
How do we guarantee safety when we deploy RL in real-world applications?

S. P. Vimal, Department of CSIS, WILP Division (vimalsp@wilp.bits-pilani.ac.in)


 **Agenda for the session**


- What do we mean by **Safety** in RL?
- Safe RL – Problem
- Some Applications (Autonomous Driving)

Source for the session & Recommended Reading:
A Review of Safe Reinforcement Learning: Methods, Theory and Applications, 2023, Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, Alois Knoll [Available from: <https://arxiv.org/pdf/2205.10330.pdf>]

S. P. Vimal, Department of CSIS, WILP Division (vimalsp@wilp.bits-pilani.ac.in)

2

 **Safety**



BIG QUESTION:
How do we guarantee **safety** when we deploy RL in real-world applications?

S. P. Vimal, Department of CSIS, WILP Division (vimalsp@wilp.bits-pilani.ac.in)

Safety

Some Definitions of Safety:

- (1) the condition of being protected from or unlikely to cause danger, risk, or injury [oxford]
- (2) being protected from harm or other dangers
- (3) controlling recognized dangers to attain an acceptable level of risk

S. P. Vimal, Department of CSIS, WILP Division (vimalsp@wilp.bits-pilani.ac.in)

Safety

2H3W problems:

- (1) **Safety Policy.** How can we perform policy optimisation to search for a safe policy?
- (2) **Safety Complexity.** How much training data is required to find a safe policy?
- (3) **Safety Applications.** What is the up to date progress of safe RL applications?
- (4) **Safety Benchmarks.** What benchmarks can we use to fairly and holistically examine safe RL performance?
- (5) **Safety Challenges.** What are the challenges faced in future safe RL research?

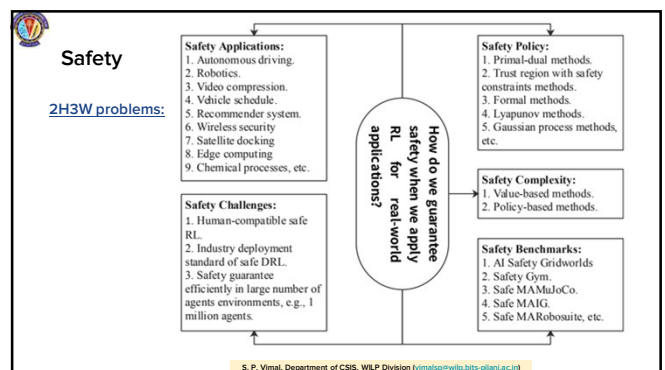
S. P. Vimal, Department of CSIS, WILP Division (vimalsp@wilp.bits-pilani.ac.in)

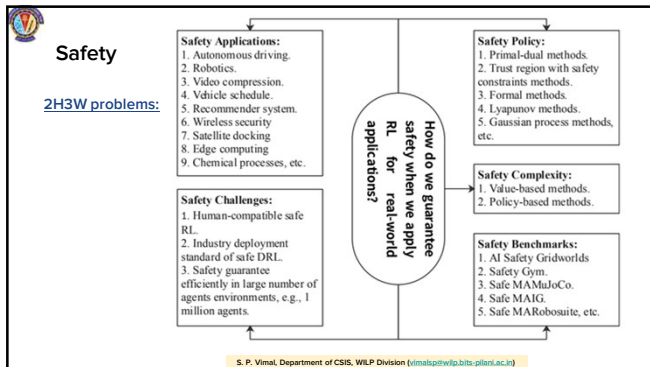
Safety

Some Definitions of Safety in the RL Sense:

- (1) humans need to label environmental states as "safe" or "unsafe," and agents are considered "safe" if "they never reach unsafe states".
- (2) agents are safe if "they act, reason, and generalize obeying human desires"

S. P. Vimal, Department of CSIS, WILP Division (vimalsp@wilp.bits-pilani.ac.in)





Defining SafeRL

MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \rho_0, \gamma)$

State Value Function $V_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s \right]$

State-Action Value Function $Q_{\pi}(s, a) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s, a_0 = a \right]$

Advantage Function $A_{\pi}(s, a) = Q_{\pi}(s, a) - V_{\pi}(s)$

Cost Function $J(\pi) = \mathbb{E}_{s \sim \rho_0(\cdot)} [V_{\pi}(s)]$

S. P. Vimal, Department of CSIS, WLP Division (vimalsp@wlp.bits-pilani.ac.in)

Defining SafeRL

MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \rho_0, \gamma)$

$\mathbb{P}(s' | s, a)$ *Transition Model*

$r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ *Reward Function*

$\rho_0(\cdot) : \mathcal{S} \rightarrow [0, 1]$ *Starting State distribution*

γ *discount factor*

S. P. Vimal, Department of CSIS, WLP Division (vimalsp@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathbb{P}, r, \rho_0, \gamma)$

Constrained

MDP with safety constraints

S. P. Vimal, Department of CSIS, WLP Division (vimalsp@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP $\mathcal{M} = (S, \mathcal{A}, \mathbb{P}, r, \rho_0, \gamma)$

Constrained

MDP with safety constraints

Cost Value function for Constraint type i

$\{(c_i, b_i)\}_{i=1}^m$

Safety Constraint Bound.

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

Now define $V_\pi^{c_i}, Q_\pi^{c_i}, A_\pi^{c_i}$ for each c_i

$V_\pi^{c_i}(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t c_i(s_t, a_t) | s_0 = s \right]$

Cost Value function for Constraint type i

$\{(c_i, b_i)\}_{i=1}^m$

Safety Constraint Bound.

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

Now define $V_\pi^{c_i}, Q_\pi^{c_i}, A_\pi^{c_i}$ for each c_i

Cost Value function for Constraint type i

$\{(c_i, b_i)\}_{i=1}^m$

Safety Constraint Bound.

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

Now define $V_\pi^{c_i}, Q_\pi^{c_i}, A_\pi^{c_i}$ for each c_i

$V_\pi^{c_i}(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t c_i(s_t, a_t) | s_0 = s \right]$

$\{(c_i, b_i)\}_{i=1}^m$

Compare!

$V_\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s \right]$

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

Now define $V_\pi^{c_i}, Q_\pi^{c_i}, A_\pi^{c_i}$ for each c_i

$$V_\pi^{c_i}(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t c_i(s_t, a_t) \mid s_0 = s \right]$$

$\{(c_i, b_i)\}_{i=1}^m$

Expected Cost function

$$C_i(\pi) = \mathbb{E}_{s \sim \rho_0(\cdot)} [V_\pi^{c_i}(s)]$$

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

Constraint Bound

$$V_\pi^{c_i}(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t c_i(s_t, a_t) \mid s_0 = s \right]$$

$$\Pi_C = \bigcap_{i=1}^m \{ \pi \in \Pi_S \text{ and } C_i(\pi) \leq b_i \}$$

Feasible Policy Set

Expected Cost function

$$C_i(\pi) = \mathbb{E}_{s \sim \rho_0(\cdot)} [V_\pi^{c_i}(s)]$$

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

$$V_\pi^{c_i}(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t c_i(s_t, a_t) \mid s_0 = s \right]$$

Expected Cost function

$$C_i(\pi) = \mathbb{E}_{s \sim \rho_0(\cdot)} [V_\pi^{c_i}(s)]$$

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

$$\Pi_C = \bigcap_{i=1}^m \{ \pi \in \Pi_S \text{ and } C_i(\pi) \leq b_i \}$$

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

$\max_{\pi \in \Pi_S} J(\pi)$, such that $c(\pi) \leq b$
 (Maximize reward performance) (Subject to Con constraint)

$\Pi_C = \cap_{i=1}^m \{ \pi \in \Pi_S \text{ and } C_i(\pi) \leq b_i \}$
 (Vector of C_i 's) (Vector of bounds)

$\pi_\star = \arg \max_{\pi \in \Pi_C} J(\pi)$
 (CMDP Optimization Function)

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

SafeRL – Categorization

Model Based SafeRL

Model-free SafeRL

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Defining SafeRL

CMDP

$\max_{\pi \in \Pi_S} J(\pi)$, such that $c(\pi) \leq b$
 (Maximize reward performance) (Subject to Con constraint)

$\Pi_C = \cap_{i=1}^m \{ \pi \in \Pi_S \text{ and } C_i(\pi) \leq b_i \}$
 (Vector of C_i 's) (Vector of bounds)

$\pi_\star = \arg \max_{\pi \in \Pi_C} J(\pi)$
 (CMDP Optimization Function)

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)

Applications | Autonomous Driving

Safe Reinforcement Learning on Autonomous Vehicles

David Ise, Alireza Nikbazi, and Kikuo Fujimura
Honda Research Institute USA
{disele, anikbazi, kfuji@hri.honda-ri.com}

Abstract—There have been numerous advances in reinforcement learning, but the typically unconstrained exploration of the learning process prevents the adoption of these methods in many safety critical applications. Recent work in safe reinforcement learning uses idealized models to achieve their guarantees, but these models do not easily accommodate the stochasticity or high-dimensionality of real world systems. We investigate how prediction provides a general and intuitive framework to constraint exploration, and show how it can be used to safely learn intersection handling behaviors on an autonomous vehicle.


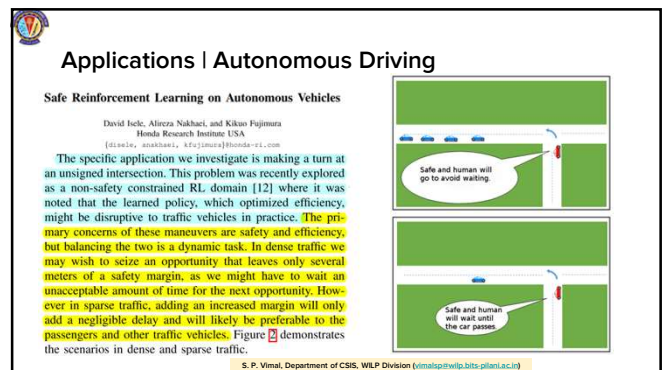
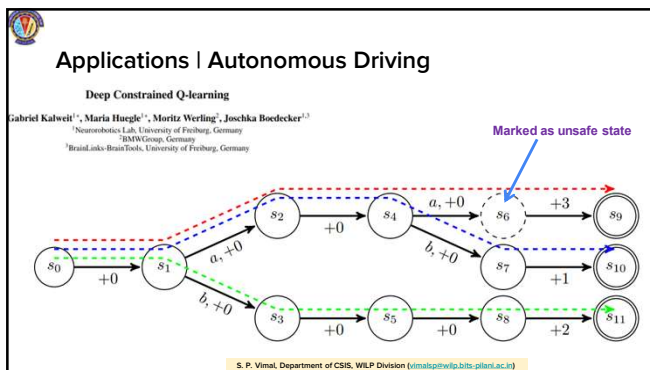
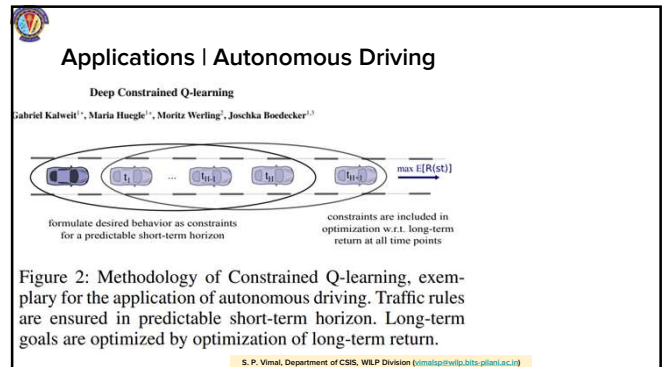
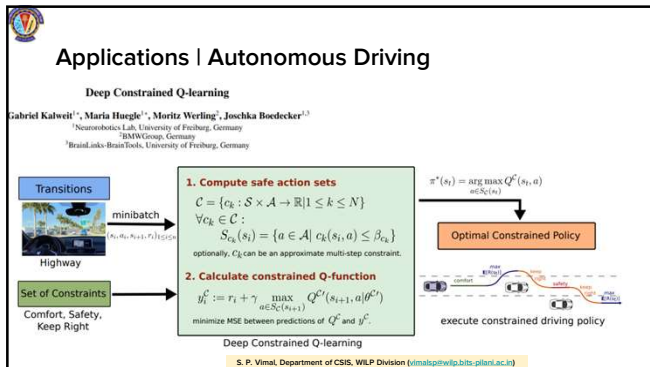


Fig. 1. An autonomous vehicle navigating an intersection. Prediction is used to shield the vehicle from making dangerous decisions, while allowing it to learn policies that are both efficient and not disruptive to other vehicles.

S. P. Vimal, Department of CSIS, WLP Division (vimal@wlp.bits-pilani.ac.in)



Applications | Autonomous Driving

Safe Reinforcement Learning on Autonomous Vehicles

David Isele, Alireza Nakhaei, and Kikuo Fujimura
Honda Research Institute USA
[disele, anakhaei, kfujimura@honda-ri.com]

<https://arxiv.org/pdf/1910.00399.pdf>

S. P. Vimal, Department of CSIS, WILP Division (vimalsp@wilp.bits-pilani.ac.in)

BITS Pilani
Pursuing the Frontiers of Knowledge

Thank you

31

Required Readings and references

1. Deep Constrained Q-learning Gabriel Kalweit, Maria Huegle, Moritz Werling, Joschka Boedecker
2. Safe Reinforcement Learning on Autonomous Vehicles David Isele, Alireza Nakhaei, and Kikuo Fujimura Honda Research Institute USA
3. A Review of Safe Reinforcement Learning: Methods, Theory and Applications, 2023, <https://arxiv.org/pdf/2205.10330.pdf>

32