Reinforcement Learning with Almost Sure Constraints

Enrique Mallada

JOHNS HOPKINS



A. Castellano H. Min J. Bazerque

Workshop on Control and Machine Learning: Challenges and Progress Oct 12, 2023

A World of Success Stories

2017 Google DeepMind's DQN



2017 AlphaZero – Chess, Shogi, Go

Boston Dynamics

2019 AlphaStar – Starcraft II



OpenAI – Rubik's Cube





Waymo





Angry Residents, Abrupt Stops: Waymo Vehicles Are Still Causing Problems in Arizona

RAY STERN | MARCH 31, 2021 | 8:26AM

GARY MARCUS BUSINESS 08.14.2019 09:00 AM

DeepMind's Losses and the Future of Artificial Intelligence

Alphabet's DeepMind unit, conqueror of Go and other games, is losing lots of money. Continued deficits could imperil investments in Al.

AARIAN MARSHALL BUSINESS 12.07.2020 04:06 PM

<u> Ilber Gives IIn on the Self-Driving Dream</u>

Can we adapt reinforcement learning algorithms to address physical systems challenges?





woman did not recognize that pedestrians jaywalk

The automated car lacked "the capability to classify an object as a pedestrian unless that object was near a crosswalk," an NTSB report said.



Challenges of RL for Physical Systems

- Physical systems must meet multiple objectives
 - Need to trade off between the different goals
 - Constrained RL allows to explore the Pareto Front [1,2]

$$\max_{\pi} \mathbb{E}_{\pi, S_0 \sim q} \left[\sum_{t=0}^{+\infty} \gamma^t R_{t+1}^{(0)} \right]$$

s.t.
$$\mathbb{E}_{\pi, S_0 \sim q} \left[\sum_{t=0}^{+\infty} \gamma^t R_{t+1}^{(i)} \right] \ge c_i, \ \forall i \in [n]$$

- Failures have a qualitatively different impact
 - Expectation constraints cannot meet safety requirements
 - Hard (almost sure) constraints can guarantee safety [3,4,5]

$$\max_{\pi} \mathbb{E}_{\pi, S_0 \sim q} \left[\sum_{t=0}^{+\infty} \gamma^t R_{t+1} \right]$$

s.t. $\mathbb{P}_{\pi, S_0 \sim q} \left[S_t \notin \mathcal{G} \right] = 1, \ \forall t \ge 0$

Zheng, You, and M, Constrained reinforcement learning via dissipative saddle flow dynamics, Asilomar 2022
 You, and M, Saddle flow dynamics: Observable certificates and separable regularization, ACC 2021
 Castellano, Min, Bazerque, M, Reinforcement Learning with Almost Sure Constraints, L4DC 2022
 Castellano, Min, Bazerque, M, Learning to Act Safely with Limited Exposure and Almost Sure Certainty, IEEE TAC, 2023
 Castellano, Min, Bazerque M, Correct-by-design Safety Critics Using Non-contractive Bellman Operators, submitted



Safety-critical Constraints in Dynamical Systems

Reachability Theory^[1-2]

- Model-based: Via Hamilton Jacobi Issacs Equations (cont. time), or iterative set updates (discrete time).
- Constraints: Provides hard/almost sure guarantees
- **Output:** Finds the maximum control invariant set (M-CIS) outside G

Control Barrier Functions (CBF)^[3-4]

- Model-based: Requires knowledge of dynamics and finding such CBF!
- **Constraints:** Provides hard/almost sure guarantees
- Output: Possibly conservative CIS

Safety Critics (SC)^[5-7]

- **Model-free:** Q-Learning-like algorithms, computes function such that $Q_{safe}(s, a) \ge \eta_{thresh} \Rightarrow$ "safety"
- **Constraints:** Provides soft/approximate guarantees, depending on discounting factor $\gamma \in (0,1)$
- **Output:** Converges to maximum CIS as $\gamma \rightarrow 1$

[1] I Mitchell, A Bayen, and C Tomlin. "A time-dependent Hamilton-Jacobi formulation of reachable sets for continuous dynamic games." IEEE TAC, 2005

[2] D Bertsekas. "Infinite time reachability of state-space regions by using feedback control." IEEE TAC, 1972

^[3] A Ames, X Xu, J Grizzle, and P Tabuada, "Control barrier function based quadratic programs for safety critical systems," IEEE TAC, 2017.

^[4] A Ames, S Coogan, M Egerstedt, G Notomista, K Sreenath, and P Tabuada. "Control barrier functions: Theory and applications" ECC, 2019

^[5] J Fisac, N Lugovoy, V Rubies-Royo, S Ghosh, and C Tomlin, "Bridging Hamilton-Jacobi safety analysis and reinforcement learning," ICRA, 2019.

^[6] K Srinivasan, B Eysenbach, S Ha, J Tan, and C Finn. "Learning to be safe: Deep RL with a safety critic." arXiv preprint arXiv:2010.14603 (2020).

^[7] B Thananjeyan, A Balakrishna, S Nair, M Luo, K Srinivasan, M Hwang, J E Gonzalez, J Ibarz, C Finn, and K Goldberg. Recovery RL: Safe reinforcement learning with learned recovery zones. IEEE Robotics and Automation Letters, 2021

Safety-critical Constraints in Dynamical Systems

Reachability Theory^[1-2]

- Model-based: Via Hamilton Jacobi Issacs Equations (cont. time), or iterative set updates (discrete time).
- Constraints: Provides hard/almost sure guarantees
- **Output:** Finds the maximum control invariant set (M-CIS) outside G

Control Barrier Functions (CBF)^[3-4]

- Model-based: Requires knowledge of dynamics and finding such CBF!
- **Constraints:** Provides hard/almost sure guarantees
- Output: Possibly conservative CIS

Safety Critics (SC)^[5-7]

- **Model-free:** Q-Learning-like algorithms, computes function such that $Q_{safe}(s, a) \ge \eta_{thresh} \Rightarrow$ "safety"
- **Constraints:** Provides soft/approximate guarantees, depending on discounting factor $\gamma \in (0,1)$
- **Output:** Converges to maximum CIS as $\gamma \rightarrow 1$

Method	Model-free	Constraint Type	Size	Control Invariant?
Reachability Theory ^[1-2]	No	Hard	Maximum CIS	Yes
Control Barrier Functions ^[3-4]	Νο	Hard	Subset of M-CIS	Yes
Safety Critics ^[5-7]	Yes	Soft/Approx.	Maximum CIS as $\gamma ightarrow 1$	No

Safety-critical Constraints in Dynamical Systems

Reachability Theory^[1-2]

- Model-based: Via Hamilton Jacobi Issacs Equations (cont. time), or iterative set updates (discrete time).
- Constraints: Provides hard/almost sure guarantees
- **Output:** Finds the maximum control invariant set (M-CIS) outside G

Control Barrier Functions (CBF)^[3-4]

- Model-based: Requires knowledge of dynamics and finding such CBF!
- **Constraints:** Provides hard/almost sure guarantees
- Output: Possibly conservative CIS

Safety Critics (SC)^[5-7]

- **Model-free:** Q-Learning-like algorithms, computes function such that $Q_{safe}(s, a) \ge \eta_{thresh} \Rightarrow$ "safety"
- **Constraints:** Provides soft/approximate guarantees, depending on discounting factor $\gamma \in (0,1)$
- **Output:** Converges to maximum CIS as $\gamma \rightarrow 1$

Method	Model-free	Constraint Type	Size	Control Invariant?
Reachability Theory ^[1-2]	No	Hard	Maximum CIS	Yes
Control Barrier Functions ^[3-4]	No	Hard	Subset of M-CIS	Yes
Safety Critics ^[5-7]	Yes	Soft/Approx.	Maximum CIS as $\gamma o 1$	No
Ours	Yes	Hard	M-CIS and Subsets	Yes



[Submitted on 18 May 2021 (v1), last revised 25 May 2021 (this version, v2)]

Learning to Act Safely with Limited Exposure and Almost Sure Certainty

Agustin Castellano, Hancheng Min, Juan Bazerque, Enrique Mallada





Agustin Castellano





Hancheng Min





Juan Bazerque



Learning for Safety-critical Sequential Decision Making



Methodology:

• Enhance RL with **logical** feedback naturally arising from constraint violations $S_t \in \mathcal{G} \Leftrightarrow D_t = 1$

Develop algorithms for learning fixed points of non-contractive operators

Recap: RL with Almost Sure Constraints

$$\max_{\pi} \mathbb{E}_{\pi, S_0 \sim q} \left[\sum_{t=0}^{+\infty} \gamma^t R_{t+1} \right]$$

s.t. $\mathbb{P}_{\pi, S_0 \sim q} \left[S_t \notin \mathcal{G} \right] = 1, \ \forall t \ge 0 \iff D_{t+1} = 0 \text{ almost surely } \forall t$



- Damage indicator $D_t \in \{0,1\}$ turns on $(D_t = 1)$ when constraints are violated
- Constraints not given a priori: Need to learn from experience!
- **Notice:** Model free → Constraint violations are inevitable

Outline

- Separation Principle for Joint Safety & Optimality
- Learning Safety with Limited Failures
- One-sided Bellman Equations for Continuous States

Formulation via hard barrier indicator

Safe RL problem:

 $V^*(s) := \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid S_0 = s \right]$

s.t.: $D_{t+1} = 0$ almost surely $\forall t$

Equivalent unconstrained formulation:

$$\max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^{t} R_{t+1} + \log[1 - D_{t+1}] \mid S_{0} = s \right]$$

$$0 \quad if \ D_{t+1} = 0$$

$$-\infty \quad if \ D_{t+1} = 1$$

Questions/Comments:

- Is this just a standard RL problem with $\tilde{R}_{t+1} = R_{t+1} + \log(1 D_{t+1})$?
- Standard MDP assumptions for Value Iteration, Bellman's Eq., Optimality Principle, etc., do not hold!
- Not to mention convergence of stochastic approximations.

Key idea: Separate the problem of safety from optimality

Hard Barrier Action-Value Functions

Consider the Q-function for a given policy π ,

$$Q^{\pi}(s,a) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \left(\gamma^{t} R_{t+1} + \log(1 - D_{t+1}) \right) \mid S_{0} = s, A_{0} = a \right]$$

and define the hard-barrier function

$$B^{\pi}(s,a) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \log(1 - D_{t+1}) \mid S_0 = s, A_0 = a \right]$$

Notes on $B^{\pi}(s, a)$:

- $B^{\pi}(s,a) \in \{0,-\infty\}$
- Summarizes safety information
 - $B^{\pi}(s, a) = 0$ iff π is safe after choosing $A_t = a$ when $S_t = s$
- It is independent of the reward process

Separation Principle

Theorem (Separation principle)

Assume rewards R_{t+1} are bounded almost surely for all t. Then for every policy π :

$$Q^{\pi}(s,a) = Q^{\pi}(s,a) + B^{\pi}(s,a)$$

In particular, for optimal π_*

$$Q^*(s, a) = Q^*(s, a) + B^*(s, a)$$

Approach: Learn feasibility (encoded in B^*) independently from optimality.

Optimal Hard Barrier Action-Value Function

Theorem (Safety Bellman Equation for B^*) Let $B^*(s, a) := \max_{\pi} B^{\pi}(s, a)$, then the following holds: $B^*(s, a) = \mathbb{E}\left[-\log(1 - D_{t+1}) + \max_{a'} B^*(S_{t+1}, a') \mid S_0 = s, A_0 = a\right]$

Understanding **B**^{*}(s, a):

 $B^*(s, a) \in \{0, -\infty\}$ summarizes safety information of the entire MDP

- $B^*(s, a) = 0$ if \exists safe π after choosing $A_t = a$ when $S_t = s$ Control Invariant
- $B^*(s, a) = -\infty$ if no safe policy exists after choosing $A_t = a$ when $S_t = s$ Unsafe



Properties of Safety Bellman Equation

Theorem (Safety Bellman Equation for B^*) Let $B^*(s,a) := \max_{\pi} B^{\pi}(s,a)$, then the following holds: $B^*(s,a) = \mathbb{E}\left[-\log(1 - D_{t+1}) + \max_{a'} B^*(S_{t+1},a') \mid S_0 = s, A_0 = a\right]$

Understanding the Solutions to the Safety Bellman Equation (SBE):

- SBE can have **multiple solutions**, including $\widetilde{B}(s, a) = -\infty$, for all pairs (s, a)
- If the function \widetilde{B} is a solution to the SBE, then:
 - The set $C \coloneqq \{s : \max_{a} \tilde{B}(s, a) = 0\}$ is a control invariant safe set
 - C is *maximal*: If $S_0 \notin C$, then S_t never reaches C for all policies π



Outline

- Separation Principle for Joint Safety & Optimality
- Learning Safety with Limited Failures
- One-sided Bellman Equations for Continuous States

Learning the barrier in finite MDPs...

Algorithm 3: barrier_update

B-function (initialized as all-zeroes); Input: (s, a, s', d)Output: Barrier-function B(s, a) $B(s, a) \leftarrow B(s, a) + \log(1 - d) + \max_{a'} B(s', a')$

...with a generative model:

Pros:

- Wraps around learning algorithms (Q-learning, SARSA)
- Use the B to trim the exploration set and avoid repeating unsafe actions

• Sample a transition (s, a, s', d) according to the MDP. Update barrier function.



Convergence in Expected Finite Time

Theorem (Safety Guarantee): Let
$$T = \min_{t} \{B^{(t)} = B^*\}$$
, then
 $\mathbb{E}T \le (L+1) \frac{|S||A|}{\mu} \left(\sum_{k=1}^{|S||A|} \frac{1}{k}\right)$

- After $T = \min_{t} \{B^{(t)} = B^*\}$, all "unsafe" (s, a)-pairs are detected
- μ : Lower bound on the non-zero transition probability

$$u = \min\{p(s', d|s, a): p(s', d|s, a) \neq 0\}$$

• L: Lag of the MDP

 $L = \max_{\substack{(s,a)\\B^*(s,a)=-\infty}} \{ \begin{array}{c} \frac{\text{Minimum}}{\text{needed to observe damage,}} \\ \text{starting from unsafe } (s,a) \end{array} \}$

Lag of the MDP: L

$$= \max_{\substack{(s,a)\\ B^*(s,a) = -\infty}} \left\{ \begin{array}{c} \frac{\text{Minimum}}{\text{minimum}} \text{number of transitions needed to} \\ \text{observe damage, starting from unsafe} (s,a) \end{array} \right\}$$



16

Sample Complexity of Safety

Theorem (Sample Complexity): With at least $1 - \delta$ probability, the algorithm learns optimal barrier function B^* after

$$(L+1)\frac{|S||A|}{\mu}\left(\sum_{k=1}^{|S||A|}\frac{1}{k}\right)\log\frac{1}{\delta}$$

iterations

- Concentration of sum of exponential random variables
- Much more sample-efficient than "learning an ϵ -optimal policy with 1δ probability" (Li et al. 2020)

$$N = \frac{|S||A|}{(1-\gamma)^{4}\varepsilon^{2}}\log^{2}\left(\frac{|S||A|}{(1-\gamma)\varepsilon\delta}\right)$$

Sample Complexity of Safety

Theorem (Sample Complexity): With at least $1 - \delta$ probability, the algorithm learns optimal barrier function B^* after

$$(L+1)\frac{|S||A|}{\mu}\left(\sum_{k=1}^{|S||A|}\frac{1}{k}\right)\log\frac{1}{\delta}$$

iterations

- Concentration of sum of exponential random variables
- If the Barrier Function is learnt first, then learning an ϵ -optimal policy takes

$$N' = \frac{|S_{safe}||A_{safe}|}{(1-\gamma)^{4}\varepsilon^{2}}\log^{2}\left(\frac{|S_{safe}||A_{safe}|}{(1-\gamma)\varepsilon\delta}\right)$$

samples (Trimming the MDP by learning the barrier)

Numerical Experiments

Goal: Reach the end of the aisle $(R_{t+1} = 10)$

Touching the wall gives $D_{t+1} = 1$, resets the episode.

Results



Actions

 s_2

 s_1

 s_3

 s_4

. . .

Why does Assured Q-learning perform much better?

If $D_{t+1} = 1 \Longrightarrow B_{\pi}(s, a) = -\infty \Longrightarrow \underline{\text{Never}}$ take action a at s again!

Takeaways:

- Adding constraints to the problem can accelerate learning
- Barrier function avoids actions that lead to further wall bumps

 s_{14}

 s_{15}

Numerical Experiments II

Setup: Rectangular grid, stepping into **holes** gives damage $D_t = 1$.

Actions $A = \{up, down, left, right\}.$

With every action, small probability to move to a random adjacent state.

Result: Barrier-learner identifies **all** the state space as unsafe.



Numerical Experiments II

Setup: Rectangular grid, stepping into **holes** gives damage $D_t = 1$.

Actions $A = \{up, down, left, right\}.$

With every action, small probability to move to a random adjacent state.

Result: Barrier-learner identifies **all** the state space as unsafe.



Outline

- Separation Principle for Joint Safety & Optimality
- Learning Safety with Limited Failures
- One-sided Bellman Equations for Continuous States

Recall: Properties of Safety Bellman Equation

Theorem (Safety Bellman Equation for B^*) Let $B^*(s, a) := \max_{\pi} B^{\pi}(s, a)$, then the following holds: $B^*(s, a) = \mathbb{E}\left[-\log(1 - D_{t+1}) + \max_{a'} B^*(S_{t+1}, a') \mid S_0 = s, A_0 = a\right]$

Understanding the Solutions to the Safety Bellman Equation (SBE):

- SBE can have **multiple solutions**, including $\widetilde{B}(s, a) = -\infty$, for all pairs (s, a)
- If the function \widetilde{B} is a solution to the SBE, then:

• The set
$$C := \left\{ s : \max_{a} \tilde{B}(s, a) = 0 \right\}$$
 is a control invariant safe set

• -C is maximal: If $S_0 \notin C$, then S_t never reaches C for all policies π





Problem: Maximal solutions can be very close to unsafe region $\mathcal{R}(\mathcal{G})$

One-Sided Safety Bellman Equation

Theorem (One-Sided Safety Bellman Equation) Let $\tilde{B}(s, a)$ be a solution of the following set of inequalities: $\tilde{B}(s,a) \leq \mathbb{E}\left[-\log(1-D_{t+1}) + \max_{a'} \tilde{B}\left(S_{t+1},a'\right)|S_0 = s, A_0 = a\right]$ The set $\mathcal{C} \coloneqq \left\{s : \max_{a} \tilde{B}(s,a) = 0\right\}$ is a control invariant safe set, not

necessarily maximal



Learning CIS Using Deep Neural Nets

Algorithm Summary

- Uses axiomatic data $(s, a, d, s') \in \mathcal{D}_{safe}$ known to be safe
- Initialize $\hat{b}^{\theta}(s, a) = 0$, where $\hat{b}(s, a) = 1 e^{B(s, a)}$ (all presumed safe)
- At each iteration, take N episodes starting from \mathcal{D}_{safe}
 - Behavioral policy: uniform safe policy

$$\pi^{\theta}(a|s) = \begin{cases} 0 & \text{if } \hat{b}^{\theta}(s,a) = 1\\ 1/\sum_{a' \in \mathcal{A}} \mathbb{1}\{\hat{b}^{\theta}(s,a') = 0\} & \text{if } \hat{b}^{\theta}(s,a) = 0 \end{cases}$$

- Train NN using SGD until fully fitting the data
- Start a new iteration, and repeat

Numerical Illustration

Control Engineer Favorite's: Inverted Pendulum



SBE = Fisac's '19 Safety Critic

Summary and future work

- Reinforcement Learning for Safety-Critical Systems
- Treat constraints separately or in parallel (Barrier Learner)
- Can characterize all feasible policies ($D_t \equiv 0$) with finite mistakes
- Requires learning using Bellman equations with non-unique solutions
- Takeaways:
 - Learning feasible policies is simpler than learning the optimal ones
 - Adding constraints makes optimal policies, easier to find
 - One-sided Safe Bellman can be used to find CISs that are not maximal

Thanks!

Related Publications:

Castellano, Min, Bazerque, M, Reinforcement Learning with Almost Sure Constraints, L4DC, 2022
 Castellano, Min, Bazerque, M, Learning to Act Safely with Limited Exposure and Almost Sure Certainty, IEEE TAC, 2023
 Castellano, Min, Bazerque M, Correct-by-design Safety Critics Using Non-contractive Bellman Operators, submitted







Hancheng Min

Enrique Mallada

mallada@jhu.edu http://mallada.ece.jhu.edu



