

Last Time

- ☐ robust safety
- ☐ zero-sum games

This time

- ☐ Computational tools

Lecture 7

EAS S'26

Ken Nakamura

Note: We will use AI notation this lecture

s = 'state' a = 'action' / 'control'

Recap

failure set $F \Rightarrow$ never enter F $\left\{ \begin{array}{l} \text{collisions} \\ \text{exceed joint limits} \\ \text{spill etc...} \end{array} \right.$

$$F := \{ s \mid l(s) < 0 \}$$

O.C. problem

$$V(s, t) = \max_{a(\cdot) \in A} \left[\min_{\tau \in [t, T]} l(\xi_s^a(\tau)) \right]$$

Hamilton-Jacobi-Bellman Variational Inequality

$$\min \{ l(s) - V(s, t), \frac{\partial V}{\partial t} + \max_{a \in A} \frac{\partial V}{\partial s} \cdot f(s, a) \} = 0$$

$$V(s, T) = l(s)$$

cont time: $\dot{s} = f(s, a)$

Hamilton-Jacobi-Bellman Eqn

$$V_t(s) = \min \{ l(s), \max_{a \in A} V_{t+1}(f(s, a)) \}$$

$$V_T = l(s)$$

disc time: $s' = f(s, a)$

Re-organized a bit...

$$V_t(s) = \min \{ l(s), \max_{a \in A} V_{t+1}(f(s,a)) \}$$

$$0 = \min \{ l(s) - V_t(s), \max_{a \in A} \underbrace{V_{t+1}(f(s,a)) - V_t(s)}_{\text{"difference in value in 1 timestep"}} \}$$

$$0 = \min \{ l(s) - V(s,t), \underbrace{\frac{\partial V}{\partial t} + \max_{u \in U} \frac{\partial V}{\partial s} \cdot f(s,u)}_{\text{total derivative wrt time}} \}$$

$$\frac{dV(s,t)}{dt} = \frac{\partial V}{\partial t} + \frac{\partial V}{\partial s} \cdot \underbrace{\frac{\partial s}{\partial t}}_{\text{dynamics}}$$

Even when $l(s)$, $f(s,a)$ are known analytically,
really hard to get analytic expression for V

\Rightarrow need computational tools to obtain V

Today:

- \hookrightarrow grid-based numerical methods
- \hookrightarrow self-supervised learning
- \hookrightarrow reinforcement learning

} all work for HJI,
reach-avoid, etc.

Design decision: How to represent V?

HW1: Represent $V(s,t)$ as a grid

Solve via numerical methods

helper_oc (MATLAB)

BEACLS (C++)

optimized_dp (Python)

hj_reachability (JAX)

$$s \in S \subset \mathbb{R}^d \quad a \in A \quad t \in [0, T]$$

↓
d-dim state

→ HW1: $[51 \times 51 \times 51]$

- discretize each dim M times

- discretize time K times

$$Q = \min \left\{ l(s) - V(s,t), \frac{\partial V}{\partial t} + \max_{a \in A} \frac{\partial V}{\partial s} \cdot f(s,a) \right\}$$

↓
approximated w/ fancy
finite-differences
on grid nodes

"essentially non-oscillatory"
methods

- at each grid point, need to compute "Hamiltonian"

$$H(s,t) = \max_a \frac{\partial V}{\partial s} \cdot f(s,a)$$

$$s = \begin{bmatrix} x \\ y \\ \theta \end{bmatrix} \quad \dot{s} = \begin{bmatrix} \cos \theta \\ \sin \theta \\ a \end{bmatrix} \quad \frac{\partial V}{\partial s} = \begin{bmatrix} p_x \\ p_y \\ p_\theta \end{bmatrix} \quad a \in [-1, 1]$$

$$\Rightarrow H(s,t) = \max_a \left[p_x \cos \theta + p_y \sin \theta + p_\theta \cdot a \right]$$

$$\Rightarrow a^* = \text{sign}[p_\theta] \Rightarrow \begin{array}{l} \text{control affine dyn} \\ \text{w/ box ctrl limits} \\ \text{is simple sign check} \end{array}$$

Challenge Q (think at home)

$$s = \begin{bmatrix} x \\ y \end{bmatrix} \quad \dot{s} = \begin{bmatrix} \cos a \\ \sin a \end{bmatrix} \quad a \in [0, 2\pi]$$

\uparrow directly ctrl θ

What is θ^* ?

\Rightarrow need to store value @ each grid pt in time

\Rightarrow need to do optimization @ each grid pt

Q: How many grid pts?

$K m^d \Rightarrow$ expensive!

Note: you can do "smarter" gridding. See:

"Scalable learning of..."

Herbert + Choi ICRA'21

Computation time + memory requirement blows up in state dim d

\hookrightarrow "Curse of dimensionality"

Need ~ 1 TB RAM to store $V(s, t)$ on 7-8 dim state spaces

In-practice Some dimensions matter "less"

Assign higher discretization to state dims that "matter more"

Leung et al. "On infusing Reachability-based..." IJRR'14

$$\text{State} = [y, \psi_R, V_{xR}, V_{yR}, r_R]$$

$$\text{Grid} = [21, 9, 9, 9, 9] \rightarrow \text{higher fidelity for } y$$

Grid-based numerical methods

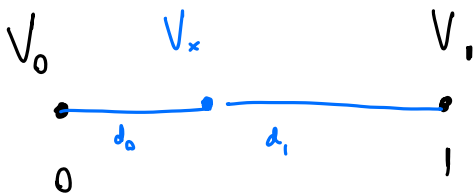
Accuracy depends on grid resolution + derivative approx

↳ can get arbitrarily accurate w/ finer grids

Computation + memory scale exponentially w/ state dim

Need to interpolate for 'off-grid' states

1-D example



$$V_x = V_0 \cdot \frac{d_1}{d_1 + d_0} + V_1 \cdot \frac{d_0}{d_1 + d_0}$$

CODE DEMO

How to scale up?

Idea: Represent $V(s, t)$ as a NN, solve w/ some algorithm

Simplest Idea: supervised learning

$$\min_{\theta} \mathbb{E}_s \left\| \underbrace{V_{g_t}(s, t)}_{\text{true val fn}} - \underbrace{V_{\theta}(s, t)}_{\text{NN approx}} \right\|$$

What's wrong w/ this?

↳ we do not have $V_{g_t}(s, t)$!

What we do know

$$V_{\theta}(s, T) = l(s)$$

$$\theta = \min \left\{ l(s) - V_{\theta}(s, t), \frac{\partial}{\partial t} V_{\theta}(s, t) + \max_a \left\langle \frac{\partial}{\partial s} V_{\theta}(s, t), f(s, a) \right\rangle \right\}$$

Idea: Just enforce these conditions!

Deepreach (Bansal + Tomlin ICRA '20)

$$L_1(s, t; \theta) = \|V_\theta(s, t) - l(s)\|$$

↳ terminal time constraint

$$L_2(s, t; \theta) = \left\| \min \left\{ l(s) - V_\theta(s, t), \frac{\partial}{\partial t} V_\theta(s, t) + \max_{a \in A} \frac{\partial}{\partial s} V_\theta(s, t) \cdot f(s, a) \right\} \right\|$$

$H(s, t)$

↳ Implicit supervision of HJ-VI

Goal: $\min_{\theta} \mathbb{E}_{(s, t)} [L_1(s, t; \theta) + \lambda L_2(s, t; \theta)]$

Q: How to sample (s, t) ?

Sampling t : "backwards in time"

→ Pretrain w/ $t=T$ → fit $V(s, T) = l(s)$

→ Uniformly sample $[t_{\text{start}}, T]$ where

$t_{\text{start}} = T \rightarrow 0$ over training

Conceptually same as grid-based methods

Sampling S: Uniform

→ assumes all states are equally important

CODE Demo

In practice

- ↳ scales to high-D (100+ dim) sys
- ↳ requires special architectures
 - ↳ sinusoidal activations
 - ↳ why not ReLU?
- ↳ still need to do inner optimization
- ↳ Weak, implicit source of supervision
 - ↳ struggle w/ large T
 - ↳ may need to modify sampling procedure over states for higher quality

So far

always need to compute $\max_a \frac{\partial V}{\partial s} \cdot f(s, a)$

→ may not have access to $\dot{s} = f(s, a)$

→ Maximization may not be easy to do

Idea 2: Learn $\operatorname{argmax}_a \frac{\partial V}{\partial s} \cdot f(s, a)$?

⇒ Off-policy Actor-Critic Reinforcement learning

↳ SAC, DDPG, TD3, DrQ etc...

⇒ discrete time: $s' = f(s, a)$ $V(s) = \min \{l(s), \max_a V(f(s, a))\}$

Reachability RL (Akametalu et al 2018)
(Fisac et al. ICRA 2019)

Jointly learn π and Q s.t.

$$Q(s, a) = \min \{l(s), \max_{a'} Q(s', a')\}$$

$$\pi(s) = \operatorname{argmax}_a Q(s, a)$$

$$V(s) = \max_a Q(s, a) \quad Q(s, a) = \min \{l(s), V(s')\}$$

But first...

Discounted Safety Bellman eq

$$Q(s, a) = (1 - \gamma) l(s) + \gamma \left[\min_{\{l(s)\}} \max_{a'} Q(s', a') \right]$$

$$\gamma \in [0, 1)$$

discount factor necessary for convergence of RL algs

① Collect and store (s, a, l, s') tuples

↳ usually using $\pi_{\psi} + \text{noise}$

② fit Q : $L(\theta) = \mathbb{E}_{(s, a, l, s')} \| Q(s, a) - y_{\text{target}} \|^2$

$$y_{\text{target}} = (1 - \gamma) l(s) + \gamma \left[\min_{\{l(s)\}} \max_{a'} Q(s', \pi_{\psi}(s')) \right]$$

③ update π_{ψ} : $L(\psi) = \mathbb{E}_s \left[-Q(s, \pi_{\psi}(s)) \right]$

④ Repeat til done

Note: at beg of training, y_{target} is inaccurate

Fisac et al. : "anneal" $\gamma: 0 \rightarrow 1$ over training

↳ Familiar?

CODE Demo

RL in practice

↳ Can often use fixed $\gamma = 0.95$ or 0.99

↳ $Q(s, a) \approx Q_\theta(s, a)$ can be unreliable, try using

$s' = f(s, a) \Rightarrow$ not always possible/desirable

$$Q(s, a) \approx \min \{ l(s), Q_\theta(s', \pi(s')) \}$$

Q: What distribution is E over?
(S, a, S')

A: distribution induced by π_γ
(and reset dist)

Follow-up Q: Where will Q_π be good?

A: States where the safety policy visits

Takeaway Learning is not a silver bullet, each method has trade-offs

	Scales?	Supervision Strength	Dynamics	Accuracy
Grid	X	✓	analytic & cont time	Grid resolution over states
Deepreach	✓	X	analytic & cont time	<u>Chosen</u> Sampling dist. over states and time
RL	✓	✓	Simulator access & discrete time	<u>Induced</u> Sampling dist over states and actions