# Last Time

- ☐ Grid-based
- ☐ Deepreach

# This time

- ☐ Reachability RL
- ☐ online updates

Project Proposal Due Feb 16

Q: "Name two ways Deepreach modernized grid-based computation"

# So far

$$0 = \min\left\{ \ell(S) - V(S,t), \frac{\partial V}{\partial t} + \max_{a} \frac{\partial V}{\partial S} \cdot f(S,a) \right\}$$

| | | | |
|---|---|---|---|
| Grid | $\max_{a} \frac{\partial V}{\partial S} \cdot f(S,a)$ | Grid nodes | backward in time |
| Deepreach | $\max_{a} \frac{\partial V}{\partial S} \cdot f(S,a)$ | MLP | time curriculum |

## Example: Computing H

$$S = \begin{bmatrix} x \\ y \\ \theta \end{bmatrix} \qquad \dot{S} = \begin{bmatrix} \cos\theta \\ \sin\theta \\ a \end{bmatrix} \qquad \frac{\partial V}{\partial S} = \begin{bmatrix} P_x \\ P_y \\ P_\theta \end{bmatrix} \qquad a \in [-1, 1]$$

$$\Rightarrow H(S,t) = \max_{a} \left[ P_x \cos\theta + P_y \sin\theta + P_\theta \cdot a \right]$$

$$\Rightarrow a^* = \text{sign}[P_\theta]$$

Control affine dyn
$\Rightarrow$ w/ box ctrl limits
is simple sign check

$$\max_{a} \frac{\partial V}{\partial s} \cdot f(s,a)$$

$\downarrow$

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

$\longrightarrow S_{t+1} = f(s_t, a_t)$ $\nearrow$ Real World rollout
$\searrow$ Simulator (eg Mujoco)
$\searrow$ World Model

① might not have $\dot{s} = f(s,a)$!

② computing $\max_{a}$ might be hard!

## Setting

$s' = f(s,a)$ <u>Goal</u> Learn $V(s)$ via experience

$V(s) = \min \{ l(s), \max_{a} V(f(s,a)) \}$ $\nearrow$ discrete time HJ

$\{ (s_t, a_t, l_t) \}_{t=0}^{T}$   $l(s_t)$ ↑

"experience" in form of trajectory rollouts



To do this, need a few modifications...

1) Time discounting

Toy Ex: $V(S) = \min_{t \in [\tau, T]} \ell(S_t)$    as $T \to \infty$, $V(S) \to -\infty$!

$T \to \infty$    $\Rightarrow$ bad

$\ell(s) > \ell(s) > \ell(s) > \ell(s)$    $\ell(s) = 0$    $> \ell(s) > \dots$

Add discount factor $\gamma \in [0, 1)$

$\hookrightarrow$ at each step, $(1-\gamma)$ chance the traj ends

$$V(S) = (1-\gamma)\, \ell(S) + \gamma \min \left\{ \ell(S), \max_a V(f(s,a)) \right\}$$

episode ends right now

w/ probability $\gamma$, we get to experience the future

Time-discount induces a "contraction"

Informally... for any initialization of $V(S)$,

repeatedly applying the update across all states

$$V(S) \leftarrow (1-\gamma)\, \ell(S) + \gamma \min \left\{ \ell(s), \max_a V(f(s,a)) \right\}$$

$V(S)$ will eventually converge to a fixed pt

$\Rightarrow$ necessary for infinite horizon RL to converge

# 2) State-action Value fn

$$V: S \to \mathbb{R}$$

"Worst case $\ell$ if gov always do $\max_a V(f(s,a))$"

$$Q: S \times A \to \mathbb{R}$$

"Worst-case $\ell$ if you first take $a$, then follow $\max_a V(f(s,a))$ for all following time"

$$Q(s,a) = (1-\gamma)\ell(s) + \gamma \min\left\{\ell(s), \max_{a' \in A} Q(s',a')\right\}$$

$$s' = f(s,a)$$

$$V(s) = \max_a Q(s,a) \qquad Q(s,a) = \min\left\{\ell(s), V(s')\right\}$$

$Q(s,a)$ allows us to learn the safety of any action from a given state $s$

# Reachability RL  (A Kametalu '18)
(Fisac '19 ICRA)

Objective: Jointly learn $Q$ and Policy $\Pi$
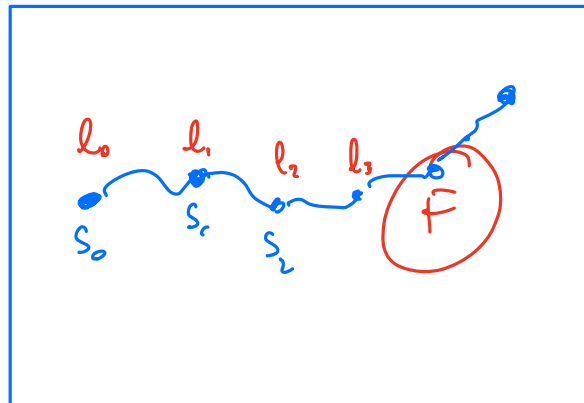
$$Q(s,a) = \cdots \text{ See above}$$

$$\Pi \approx \operatorname*{argmax}_a Q(s,a)$$

① Initialize $s_0$
"env. reset()"
↓
reset often times a
uniform distribution over states,
or near unsafe states

② take actions in environment
↳ store $(s, a, \ell, s')$ into
a Replay Buffer



③ Fit $Q$: Sample $(s, a, \ell, s')$ from Buffer

$$Q_\theta \quad L(\theta) = \mathbb{E}_{(s,a,\ell,s')} \| Q_\theta(s,a) - y_{\text{target}} \|^2$$

$$y_{target} = (1-\gamma)\,\ell(S) + \gamma\left[\min\{\ell(S), Q(S', \pi(S'))\right]$$

current policy $\pi$

don't let gradients flow through $y_{target}$

④ Improve $\pi_\psi$

$$L(\psi) = \underset{S}{\mathbb{E}}\left[-Q_Q(S, \pi_\psi(S))\right]$$

Minimize $-Q(S,a)$

$\Rightarrow$ Maximize $Q(S,a)$

⑤ Repeat until done

Note: $y_{target}$ is inaccurate estimate at the beginning of training

$$y_{target} = (1-\gamma)\,\ell(S) + \gamma\left[\min\{\ell(S), Q(S', \pi(S'))\right]$$

Current state, accurate

estimate of future
↳ bad at first!

Idea (Fisac '19): anneal $\gamma$ from low to high

e.g. $\gamma = 0.8 \rightarrow 0.9999$

$\downarrow$

Can interpret this as doing
backward-in-time computation!

---

# Robust ver. ISAACS (Hsu et al '22)

actor $\pi_{safe}(s)$, disturbance $\pi_d(s)$, $Q(s, a, d)$
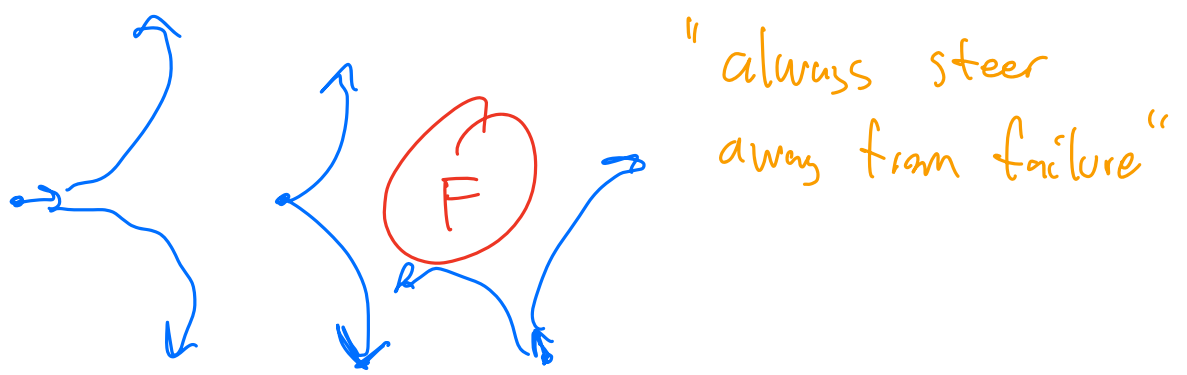
Sample $(s, a, d, l, s')$ from replay buffer

$\pi_{safe}$ loss: $-Q(s, \pi_{safe}, d)$

$\pi_d$ loss: $Q(s, a, \pi_d)$ $\longrightarrow$ update $\pi_d$ faster
to give information adv

Critic loss $\|Q(s, a, d) - y_{tar}\|$

---

Observation:

Replay Buffer is filled primarily w/
safety-preserving actions

"always steer away from failure"

$\pi_{nom} : S \to A$ is different from $\pi_{safe}$

$Q(s, \underbrace{\pi_{nom}(s)}_{\uparrow})$ might be bad / inaccurate

OOD for $Q$

Alternative:

first $s' = f(s, a)$

$\quad \hookrightarrow Q(s, a) \approx \min \{ \ell(s), \overset{\text{accurate}}{Q(s', \underbrace{\pi_{safe}(s')}_{}) } \}$

In-D for $Q$

<u>Takeaway</u> Learning is <u>not</u> a silver bullet, each method has trade-offs

| | Scales? | Supervision Stregth | 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 |