Motivation

- Measuring an unknown quantity over a Markov Chain is a typical problem.
- Access to **safe** state-action space

Goal: Use correlation to learn about **unsafe** state-action space



Figure 1. A robot measuring gas concentration at an industrial site. The gas concentration is correlated between pipes and compressors, respectively. Part of the site is restricted.

- Example in Figure 1: State-of-the-art methods do not consider correlation between unsafe and safe pipes.
- Resulting policy visits compressors more frequently than pipes (more safe compressors than pipes)
- This is sub-optimal for estimating gas concentration over the whole site.

Learning Problem

- Known and deterministic MDP $M = (S, A, P, H, f, s_0)$ with episode length H
- Unknown quantity f, part of a RKHS $f \in \mathcal{H}_k$ with known kernel
- $k((s, a), (s', a')) = \Phi(s, a)^T \Phi(s', a')$ and bound $\|f\|_{\mathcal{H}_k} \leq \frac{1}{\lambda}$
- Fixed and known unsafe set $\mathcal{B}_{unsafe} \subset \mathcal{S} \times \mathcal{A}$.
- Noisy observations of unknown quantity $y = f^T \Phi(s, a) + \epsilon$ where: $\epsilon \sim \mathcal{N}(0, \sigma^2)$
- **Fixed budget** of T trajectories of length H

Problem: Estimate a linear functional of the unknown function (potentially in an unsafe region).

$$\mathsf{find} \quad \mathbf{C}\hat{f}_{\mathcal{T}} \quad \mathsf{where:} \quad \hat{f}_{\mathcal{T}} = \arg\min_{f \in \mathcal{H}_k} \sum_{i=1}^{\mathcal{T}} \sum_{(s,a) \in \tau_i} (f^{\mathcal{T}} \Phi(s,a) - y_{i,s,a,})^2 + \lambda \|f\|_{\mathcal{H}_k}$$

As typical in Experimental Design, we minimize the second moment of the estimation error:

Safety:

- Note:

Reweighting-Based Approach (Ours):



Optimization

- Thus, we can use the machinery of Convex RL [1, 3].

[1] Elad Hazan, Sham M. Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum entropy exploration, 2019.

- [2] Mojmir Mutny, Tadeusz Janik, and Andreas Krause. Active exploration via experiment design in markov chains. 2023.
- [3] Tom Zahavy, Brendan O'Donoghue, Guillaume Desjardins, and Satinder Singh. Reward is enough for convex mdps, 2022.

Safe Active Exploration

Iason Chalas Gian Hess Alexander Spiridonov

Method

$$\mathbf{E}_{T} = \mathbb{E}_{\epsilon}[\mathbf{C}(f - \hat{f}_{T})(f - \hat{f}_{T})^{T}\mathbf{C}^{T}]$$

• We define a **safe** state-action distribution polytope:

$$egin{aligned} \mathcal{D}^{\mathit{safe}} &:= \{d| d(s,a) \geq 0, \sum_{s,a} d(s,a) = 1, \sum_{a} d(s',a) = \sum_{s,a} d(s,a) P(s'|s,a), \ &orall (s,a) \in \mathcal{B}_{\mathit{unsafe}}.d(s,a) = 0 \} \end{aligned}$$

Current (Naive) Approach:

• Mutný et al. [2] derive convex upper bound for $\mathbf{E}_{\mathcal{T}}$ • Their objective, defined over \mathcal{D}^{safe} is the following:

$$\min_{d_\pi \in \mathcal{D}^{safe}} U(d_\pi) \quad ext{with:} \quad U(d_\pi) = s \left(\sum_{(s,a) \in \mathcal{S} imes \mathcal{A}} rac{H d_\pi(s,a)}{\sigma_{s,a}^2} \Phi(s,a) \Phi(s,a)^T + rac{\lambda \mathbf{I}}{T}
ight)$$

Where $s(\cdot)$ is a convex scalarization function, such as $-\log(\det(\cdot))$

• All the features $\Phi(s, a)$ from the unsafe region cancel, since $d_{\pi}(s,a)=0 \quad orall (s,a)\in \mathcal{B}_{\mathit{unsafe}}.$

• Thus, we lose the correlation information to the unsafe region.

• Introduce a reweighting $w : S \times A \rightarrow \mathbb{R}$ to the objective

$$\min_{e \mathcal{D}^{safe}} U(d_{\pi}) \quad ext{with:} \quad U(d_{\pi}) = s \left(\sum_{(s,a) \in \mathcal{S} imes \mathcal{A}} rac{H d_{\pi}(s,a) w(s,a)}{\sigma_{s,a}^2} \Phi(s,a) \Phi(s,a)^T + rac{\lambda \mathbf{I}}{T}
ight)$$

The weights are defined as follows:

$$w(s,a) = \sum_{(s',a')\in\mathcal{B}_{unsafe}} \left|rac{\Phi(s,a)^T\Phi(s',a')}{\sqrt{\Phi(s,a)^T\Phi(s,a)}\sqrt{\Phi(s',a')^T\Phi(s',a')}}
ight|+1$$

• Weights transport information from unsafe to safe region.

• High weight \Rightarrow High correlation to unsafe region \Rightarrow Higher visitation

- \mathcal{D}^{safe} preserves convexity.
- $U(d_{\pi})$ is also convex

• We use the adaptive optimization scheme introduced by Mutný et al. [2].

- Invariant correlation structure
- Unsafe state space in grey, empirical state visitation heatmap in red
- Comparison of mean squared estimation error over the whole domain

ŢŢ Ţ	<u>لي</u> م	<u>ل</u> وہ میں جب	<u>لي</u> م	ران م	<u>ل</u> وں میں جب	<u>ر</u> ائی میں چینے	<u>لي</u> م	<u>ل</u> وں میں جب	
	<u>th</u>	$\overline{\mathbb{Q}}_{\mathcal{O}}^{\frac{1}{2}}$	Ţ Ţ	$\overline{\mathbb{Q}}_{\mathcal{C}}^{-1}$		<u>تلہ "</u>	<u>tr</u>	$\overline{\mathbb{Q}}_{\mathcal{O}}^{\frac{1}{2}}$	
	ŢŢ Ţ	<u>tr</u>	<u>tr</u>	$\overline{\psi}_{\mathcal{O}}^{-1}$			ڗ ڷ ڛؖۻ	ڛٛ	¢ ¢
	<u>Ţ</u>	<u>i</u>							
	ŢŢ Ţ	<u>i</u>						<u>i</u>	
	<u>∲</u> ~	<u>∲</u> ~	₩ [™]				₩ [™]	<u>∲</u> ~	

(a) **unweighted** objective

- (a) Unweighted objective produces policy with more trajectories from below
- (b) Weighted objective produces policy with more trajectories from top
- (c) Mean squared estimation error over executed trajectories drops faster and more consistently with reweighting

- \Rightarrow Better estimate for unsafe region and overall lower MSE
- Future experiments for more general correlation structure

FINZURICH

Experiments

• Synthetic grid world with three different features (pictograms)

ڛؖ	ڗ ڷ ڛؖ	<u>ᠳ</u>	<u>ᠳ</u>	<u>ᠳ</u>	<u>مَ</u> رَّةً م	<u>ਜ਼ੵ</u> ~	<u>ᠳ</u>	<u>₽₽°</u> ₽	<u>ش</u> م
$\overline{\psi}_{\mathcal{O}}^{-}$	<u>م</u> گ	<u> </u>	<u>مَنْ *</u>	<u> </u>		$\overline{\psi}_{\mathcal{O}}^{\frac{n}{2}}$	<u>م</u> گ	$\overline{\psi}_{\mathcal{O}}^{-}$	<u><u></u></u>
	Ţ Ţ	<u>f</u>	لل لل				لله م	Ţ Ţ	
	ڛؖ								
	$\overline{\psi}_{\mathcal{O}}$	<u>i</u>							
	Ţ Ţ	Ţ Ţ	ŢŢ Ţ	Ţ Ţ	₩ 	Ţ Ţ	ŢŢ Ţ	₩ 	

(b) weighted objective



Figure 2

Discussion

• Reweighting incentivizes visitation of the safe city state (highly correlated to unsafe region) \Rightarrow More trajectories from the top through safe city state