Overview

- by tracking the bounds with a slowly changing policy.

Motivation for Directed Exploration

Common approaches to exploration like optimistic initialization are not always viable. Hence, we would like a mechanism for directed exploration in reinforcement learning. For instance if we have access to uncertainty, $\hat{U}(s, a)$, around mean estimates $\hat{Q}(s, a)$, action selection can be greedy w.r.t. $\hat{Q}(S_t, a) + \hat{U}(S_t, a)$, which provides a highconfidence upper-bound for the best possible action in the state S_t .

Let $\tilde{Q}_t = \hat{Q}_t + \hat{U}_t$, and let π_t be the policy induced by greedy action selection on Q_t . Then, this process of action selec-Further, if we assume $v_t \sim \mathcal{N}(0, \sigma^2)$, and $\bar{\mathbf{z}}_T \stackrel{\text{def}}{=} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{z}_t$, we show tion converges to a policy that is optimal under a defined that the following holds with probability at least 1 - p: density:

$$Q^* = \operatorname*{argmax}_{Q \in Q} \int_{\mathcal{S} \times \mathcal{A}} d(s, a) Q(s, a) ds da$$

This requires three key assumptions:

- Stochastic Optimism: After some point T > 0, for all $t \geq T, \mathbb{E}[Q_t(S,A)] \geq \mathbb{E}[Q^*(S,A)].$
- Shrinking Confidence Interval Radius: $\mathbb{E}[\hat{U}_t(S, A)] \leq f(t)$ for some non-negative function f with $f(t) \rightarrow 0$.
- The bounds derived are for a stationary policy. But during control, Convergent Action Values: $|\mathbb{E}[\hat{Q}_t(S,A) - Q^{\pi_t}(S,A)]| \leq g(t)$ the policy is slowly changing, and therefore, we slowly track these for some non-negative function g with $g(t) \rightarrow 0$. upper-bounds resulting in:

Optimistic Values Theorem states that under these 3 assumptions:

$$\operatorname{Regret}(T) \stackrel{\text{\tiny def}}{=} \sum_{t=1}^{T} \mathbb{E}[Q^*(S,A)] - \mathbb{E}[Q^{\pi_t}(S,A)] \le \sum_{t=1}^{T} f(t) + g$$

Results: Comparing State-of-the-Art Exploration Methods

We compare UCLS against algorithms that use other approaches to estimate confidence intervals:

- DGPQ using GPs. [1]
- LSPI-Rmax using a measure of *knownness*. [2]
- RLSVI using Bayesian Linear Regression. [3]
- UCBootstrap using bootstrapped confidence intervals. [4]

Raksha Kumaraswamy¹, Matthew Schlegel¹, Adam White^{1,2}, Martha White¹

Context-Dependent Upper-Confidence Bounds for Directed Exploration

> We propose an incremental, model-free exploration algorithm with fast-converging upper-confidence bounds, called UCLS. > We derive confidence intervals around action-values for LSTD, and use it to provide a directed exploration signal during control

Confidence Interval Bounds for Policy Evaluation

Given the noise w.r.t. the optimal estimator \mathbf{w}^* ,

$$\mathbf{v}_{t+1} = (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top \mathbf{w}^* + \mathbf{v}_t$$

and a finite set of samples, T, with $\bar{\nu}_T \stackrel{\text{\tiny def}}{=} \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{z}_t v_t$, we show that the following holds with probability at least 1 - p:

$$\mathbf{x}^{\top}\mathbf{w}^{*} \leq \mathbf{x}^{\top}\mathbf{w}_{T} + \sqrt{\frac{p+1}{p}}\sqrt{\mathbf{x}^{\top}\mathbb{E}[\mathbf{A}_{T}^{+}\bar{\boldsymbol{\nu}}_{T}\bar{\boldsymbol{\nu}}_{T}^{\top}\mathbf{A}_{T}^{+\top}]\mathbf{x}} + O\left(\mathbb{E}[(\mathbf{x}^{\top}\boldsymbol{\epsilon}_{T}^{*})^{2}]\right)$$
(1)

$$\mathbf{x}^{\top} \mathbf{w}^{*} \leq \mathbf{x}^{\top} \mathbf{w}_{T} + \sigma \sqrt{\frac{p+1}{p}} \sqrt{\mathbf{x}^{\top} \mathbb{E}[\mathbf{A}_{T}^{+} \bar{\mathbf{z}}_{T} \bar{\mathbf{z}}_{T}^{\top} \mathbf{A}_{T}^{+\top}] \mathbf{x}} + O\left(\mathbb{E}[(\mathbf{x}^{\top} \boldsymbol{\epsilon}_{T}^{*})^{2}]\right)$$
(2)

Control with Confidence Interval Bounds

- Upper-Confidence Least Squares (UCLS) for bound in Equation
- Global Variance-Upper Confidence Bound (GV-UCB) for bound g(t)in Equation (2).



Results: Advantage of Contextual Variance

This study contrasts the advantage of contextual variance estimates (UCLS) over global variance estimates (GV-UCB).



UCLS-Linear: An Effective Linear Complexity Variant



- Pros: (1) lower computational complexity, (2) may be more amenable to changing representations.
- Cons: (1) while the performance of the two algorithms is comparable, hyper parameters that need to be tuned.

Conclusion & Future Work

- Context-based exploration is a promising direction for designing this motivation within the LSTD learning framework.
- and changing representations, (2) other approaches to promote context-based exploration at a grounded or abstract level.

References

- *Conference on Machine Learning.* 2014.
- Autonomous Agents and Multiagent Systems. 2009.
- *Conference on Machine Learning.* 2016.
- Advances in Neural Information Processing Systems. 2010.









UCLS-L experiences more regret (e.g. River Swim, Puddle World), (2) two

sample-efficient learning algorithms. UCLS is a principled application of Further lines of research include: (1) UCLS/UCLS-L with other stationary

R Grande, T Walsh, and J How. "Sample Efficient Reinforcement Learning with Gaussian Processes". In: International L Li, ML Littman, and CR Mansley. "Online exploration in least-squares policy iteration". In: International Conference on I Osband, B Van Roy, and Z Wen. "Generalization and Exploration via Randomized Value Functions.". In: International M White and A White. "Interval estimation for reinforcement-learning algorithms in continuous-state domains". In:

