EXPLAINING CONTINUAL REINFORCEMENT LEARNING AGENTSWITH SHAPLEY VALUESDaniel Beechey and Özgür ŞimşekMath
RL LABArt-ai

Why explain decision-making?

- Reinforcement learning agents typically learn to act but not to explain themselves.
- This hinders deployment in settings where accountability and trust are essential.

How to explain decision-making

- Prior work [1, 2] attributes decisions to features of an agent's observations using Shapley values [3], a theory-driven approach to fairly assigning credit.
- But computing Shapley values exactly is infeasible in real-world settings.

OUR CONTRIBUTION

FastSVERL: A scalable method for explaining decision-making by attributing actions to features of an agent's observations.



0	1	2		0	1	2		Influ
0	1			0	1			ence
0	1	1	1	0	1	1	1	Low

You're playing Minesweeper—what's your next move?

FastSVERL handles off-policy data and adapts to changing behaviour, enabling Shapley-based interpretability in practical reinforcement learning settings.

How it works (if y	ou're interested)
What are Shapley-based explanations?	Approximating the characteristic function
We explain an agent's decision by attributing how each feature influences the probability of taking an action, $\pi(s,a)$.	We approximate the characteristic function $\tilde{\pi}_s^a(C)$ with a model $\hat{\pi}(s, a C; \beta)$ trained to minimise prediction error:
To do so, we consider how the action probability changes when different features are known or unknown. This is captured by a characteristic function $\tilde{\pi}_s^a(\mathcal{C})$, which measures the expected probability of action a when only features in $\mathcal{C} \subseteq \mathcal{F}$ are known:	$\mathcal{L}(eta) = \mathop{\mathbb{E}}\limits_{p^{\pi}(s)} \mathop{\mathbb{E}}\limits_{\mathrm{Unif}(a)} \mathop{\mathbb{E}}\limits_{\mathrm{Unif}(\mathcal{C})} \left \pi(s,a) - \hat{\pi}(s,a \mathcal{C};eta) ight ^2$

 $\pi(a) = \pi[(a) + a^{2} - a^{2} + a^{2}$

$$ilde{\pi}^a_s(\mathcal{C}) = \mathbb{E}\left[\pi(S,a) \mid S^{\mathsf{c}} = s^{\mathsf{c}}
ight] = \sum_{s \in \mathcal{S}^+} p^{\pi}(s \mid s^{\mathsf{c}}) \, \pi(s,a)$$

Shapley values assign credit to each feature based on its average marginal contribution across all feature subsets:

$$\phi^i(ilde{\pi}^a_s) = \sum_{\mathcal{C}\subseteq\mathcal{F}\setminus\{i\}} rac{|\mathcal{C}|!\cdot(|\mathcal{F}|-|\mathcal{C}|-1)!}{|\mathcal{F}|!} [ilde{\pi}^a_s(\mathcal{C}\cup\{i\})- ilde{\pi}^a_s(\mathcal{C})]$$

These values uniquely satisfy axioms formalising fair credit assignment.

But exact computation is infeasible in complex settings: the total cost per explanation is $\mathcal{O}(2^{|\mathcal{F}|} \cdot |\mathcal{S}|)$, $2^{|\mathcal{F}|}$ expectations over the state space \mathcal{S} .



Email:

Website:

Approximating Shapley values

We approximate the Shapley value summation with a second model $\hat{\phi}(s, a; \theta) : S \times A \to \mathbb{R}^{|\mathcal{F}|}$, trained to minimise a least-squares objective:

$$\mathcal{L}(heta) = \mathop{\mathbb{E}}\limits_{p^{\pi}(s)} \mathop{\mathbb{E}}\limits_{\mathrm{Unif}(a)} \mathop{\mathbb{E}}\limits_{p(\mathcal{C})} \Bigl| ilde{\pi}^a_s(\mathcal{C}) - ilde{\pi}^a_s(\emptyset) - \sum_{i \in \mathcal{C}} \hat{\phi}^i(s,a; heta) \Bigr|^2$$

where
$$p(\mathcal{C}) \propto rac{n-1}{inom{n}{|\mathcal{C}|} \cdot |\mathcal{C}| \cdot (n-|\mathcal{C}|)}$$

These models amortise the cost of Shapley value approximation across all states and actions.



djeb20@bath.ac.ul

djeb20.github.io

 Beechey, D., Smith, T.M. and Şimşek, Ö., 2023, July. Explaining reinforcement learning with Shapley values. In International Conference on Machine Learning (pp. 2003-2014). PMLR.
 Beechey, D., Smith, T. and Şimşek, Ö., 2025. A Theoretical Framework for Explaining Reinforcement Learning with Shapley Values. arXiv preprint arXiv:2505.07797.
 Lloyd S Shapley. A value for n-person games. Contributions to the Theory of Games, 2(28):307–317, 1953.