

Adversarial Learning of Distributional Reinforcement Learning

Yang Sui¹ Yukun Huang¹ Hongtu Zhu² Fan Zhou¹

¹Shanghai University of Finance and Economics

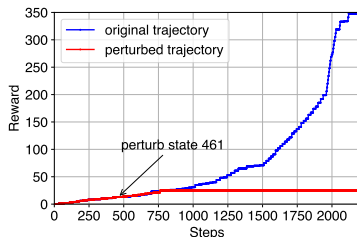
²The University of North Carolina at Chapel Hill

Motivation: Ride-Sharing

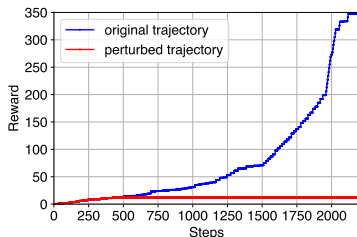


- The order-dispatching policy is trained **offline** which is then applied to real world
- Although offline and online environments may be **almost the same**, these policies can perform poorly in real world

Motivation: Empirical Studies



perturb state



perturb parameter

- Small perturbation imposed on certain **state** can significantly change the trajectory afterwards
- Similar result after perturbing some **parameter** in the policy network

Our Contributions

- Question: **How to quantify the effects of these variations?**
- Prior Work: **distribution shift between online and offline data**
- Our Contributions:
 - analyse why a trained policy fail when applied to **new but similar** environment
 - construct adversarial learning to evaluate **sensitivity** of RL components
 - our method is efficient for detecting DRL, **theoretically and empirically**

Method: Perturbation Manifold

- **DRL Setting:** Given state s , action a and network parameter θ , the distribution of z is $P(z|s, a, \theta)$.

The **perturbed model** is $P(z|s, a, \theta, \omega)$ after imposing a perturbation $\omega = (\omega_1, \dots, \omega_p)^T$ on either s or θ

- **Perturbation Manifold:** $M = \{P(z|s, a, \theta, \omega) : \omega \in \Omega\}$. The vector space of M at ω is spanned by p functions $\{\partial_i \ell(\omega|z, s, a, \theta)\}_{i=1}^p$, where $\partial_i = \partial/\partial\omega_i$ and $\ell(\omega|z, s, a, \theta) = \log P(z|s, a, \theta, \omega)$ [1, 2]
- **Metric Tensor:** $G(\omega) \in \mathbb{R}^{p \times p}$ of ω is defined as

$$g_{ij}(\omega) = \mathbf{E}_\omega [\partial_i \ell(\omega|z, s, a, \theta) \partial_j \ell(\omega|z, s, a, \theta)]$$

Influence measure

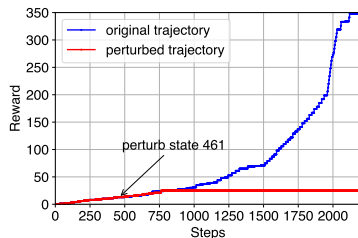
Let $f(\boldsymbol{\omega}) : \mathbb{R}^p \rightarrow \mathbb{R}^1$ be the objective function, we can then define the first-order local influence measure (**FI**) of $f(\boldsymbol{\omega})$ at $\boldsymbol{\omega}_0$ as

$$\mathbf{FI}_{\boldsymbol{\omega}}(\boldsymbol{\omega}_0) = \nabla_{f(\boldsymbol{\omega}_0)}^T \mathbf{G}^{-1}(\boldsymbol{\omega}_0) \nabla_{f(\boldsymbol{\omega}_0)}$$

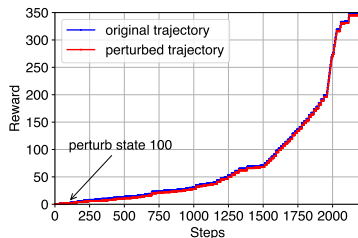
- High **FI** represents that $f(\boldsymbol{\omega})$ is more **vulnerable** to $\boldsymbol{\omega}$
- **FI** possesses an **intrinsic** property that is free of the constraints imposed by $\boldsymbol{\omega}$
- **FI** is invariant to any **reparameterization** of $\boldsymbol{\omega}$

Task (i): Detection of fragile states

Comparison of trajectories for states (Q -values as objective function)



state 461 with **high FI**:
significant difference

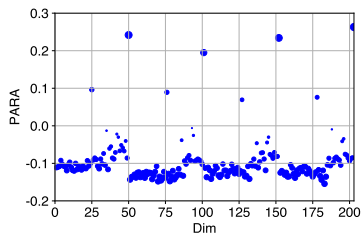


state 100 with **low FI**:
negligible difference

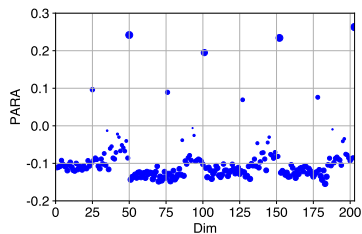
FI can detect **potentially vulnerable states!**

Task (ii): Adversarial learning of policy network

Parameters in the DENSE2 BIAS layer and their corresponding **FIs**



parameters

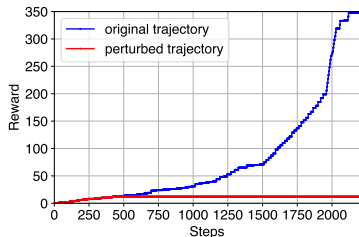


FIs

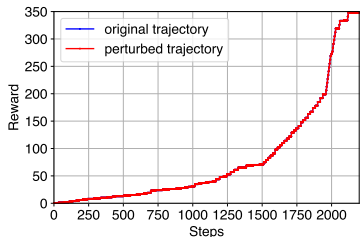
Only 8 **positive** parameters \iff 8 parameters with **high FI** detected

Task (ii): Adversarial learning of policy network

Comparison of trajectories for parameters (Q -values as objective function)



204th parameter with **high FI** :
significant difference



152nd parameter with **low FI**:
no difference

FI can pinpoint **sensitive policy parameters!**

- [1] H. Zhu, J. G. Ibrahim, S. Lee, and H. Zhang. Perturbation selection and influence measures in local influence analysis. *The Annals of Statistics*, 35(6):2565–2588, 2007.
- [2] H. Zhu, J. G. Ibrahim, and N. Tang. Bayesian influence analysis: a geometric approach. *Biometrika*, 98(2):307–323, 2011.