Asymptotically Exact Error Characterization of Offline Policy Evaluation with Misspecified Linear Models

Kohei Miyaguchi, IBM Research

miyaguchi@ibm.com



Previous studies & our contribution



- Example the second sec fundamental OPE method with linear function approximation.
- However, it has been not well understood when and how exactly linear DM works.

- Our contribution is the exact asymptotic analysis of linear <u>DM</u>.
- \rightarrow As a result, linear DM is shown to be **more robust** than previously known: *Either* Q- *or* W-realizability is sufficient to solve OPE.



• OPE = estimating policy value with other policies'

> pre-deployment performance analysis & hyperparameter tuning.

Main result

Algorithm: Linear DM

Input: data D_n , policy π , initial density p_0 , discount factor γ , feature ϕ . **Output**: policy value estimate $\hat{J}(\pi)$

> 1. $\hat{Q} \coloneqq \text{LSTD}Q(D_n, \pi, \gamma, \phi)$ 2. $\hat{J}(\pi) \coloneqq \hat{Q}(p_0, \pi)$

Theorem (asymptotic error)

If ϕ is <u>compatible</u> (and standard assumptions of OPE hold), then:

$$\hat{J}(\pi) - J(\pi) = -\frac{1}{1 - \gamma} \mathbb{E}^{q} \left[\mathcal{R} \right]$$
Error of linear DM

Definition (residual functions)

- W-residual is given by **the** projection residual of the marginal density ratio of Episode(π) to Episode(q).
- Q-residual is given by **the** projection residual of the **Q-function** of π .



Implication (double robustness)

- The residual functions being zero are equivalent to Wand Q-realizabilities.
- Linear DM is doubly robust and could be consistent even if LSTDQ is not, i.e., $Q^{\pi} \notin \operatorname{span}(\phi).$



Definition (compatibility)

 ϕ is <u>compatible</u> iff LSTDQ (D_n, π, γ, ϕ) converges as $n \to \infty$, i.e.,

$$\det \mathbb{E}^{q}[\phi(\phi - \gamma P^{\pi}\phi)^{\top}] \neq 0$$

Compatibility is **always satisfied** with onehot features ϕ .

Compatibility is **statistically testable** via concentration of determinant.

Compatibility is **non-vacuous**: It eliminates the hard instance of [Amortila+2020].



Proposition (consistency *without* realizability)

Linear DM with the homogeneous tilecoding feature ϕ is consistent if

Episode(q) is sufficiently strongly mixing and near-uniform.

2. The number of tiles grows steadily, e.g.,

$$X = \Theta(\sqrt{n})$$



• Double robustness of nonlinear estimators (e.g., MQL, MWL, FQE) Nonparametric consistency with non-onehot features (e.g., NTK)