

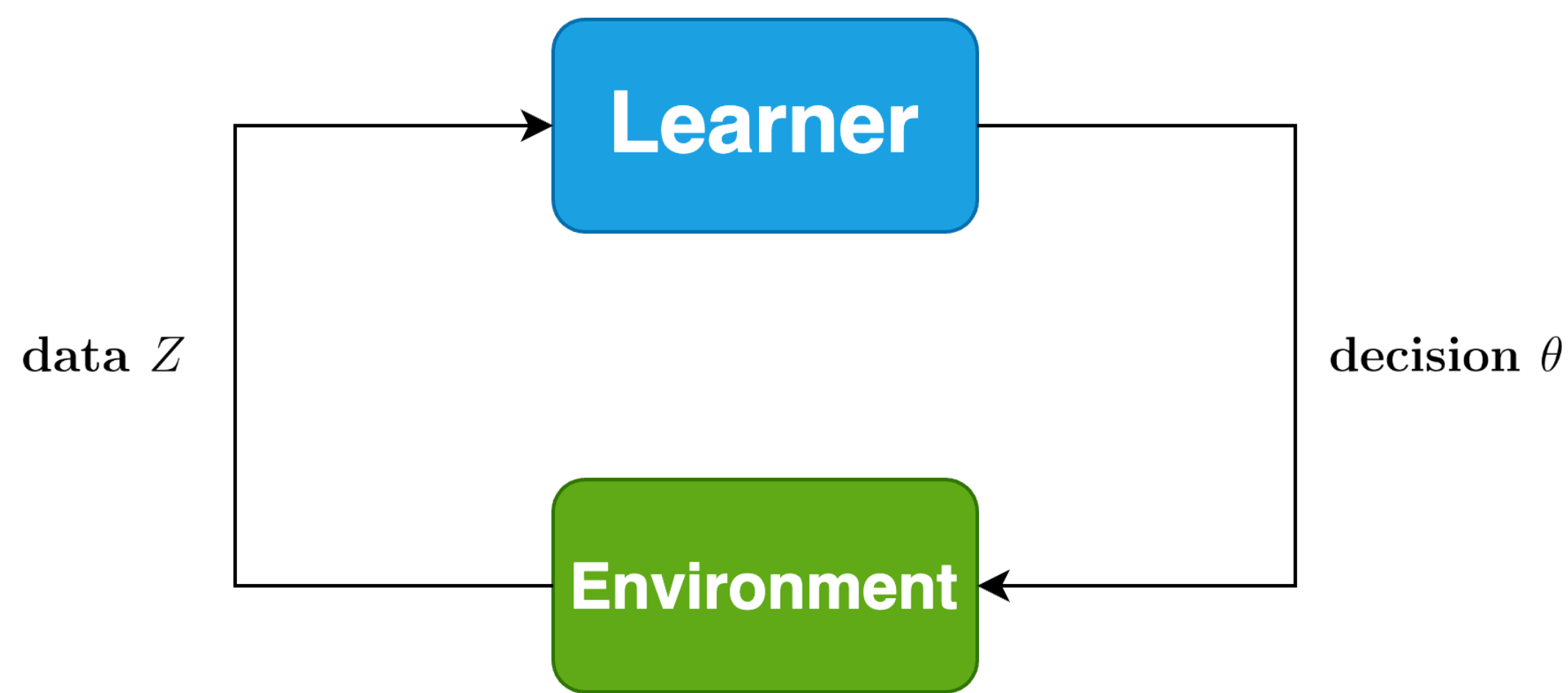
Two-timescale Derivative Free Optimization for Performative Prediction with Markovian Data

Haitong Liu, Dept. of CS, ETH Zurich, Qiang Li, Hoi-To Wai Dept. of SEEM, CUHK



Performative Prediction

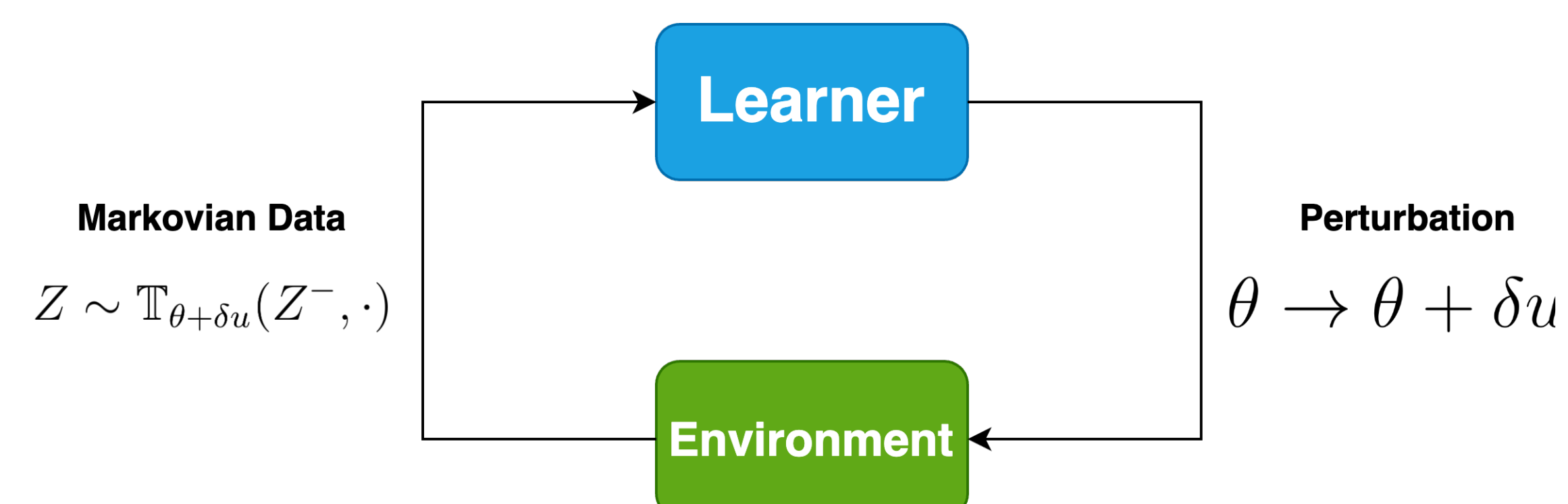
- ◇ **Performative Prediction:** data distribution depends on decision variables.
- ◇ **Motivating Examples:** loan classification, pricing, ride sharing.



- ◇ **Goal:** minimize the performative risk, $\min_{\theta} \mathcal{L}(\theta) := \mathbb{E}_{Z \sim \Pi(\theta)}[\ell(\theta; Z)] \rightarrow \text{ncvx}$
- * Evaluate $\nabla \mathcal{L}(\theta)$ needs **known** $\Pi_{\theta}(\cdot)$: $\mathbb{E}_{Z \sim \Pi_{\theta}}[\nabla \ell(\theta; Z) + \ell(\theta; Z) \nabla \log \Pi_{\theta}(Z)]$.

Zeroth Order Oracle & Markovian Data

- ◇ **One-point gradient estimator** $g_{\delta}(\cdot)$. Absence of prior knowledge on $\Pi_{\theta} \rightarrow$ **deploy and observe** $\ell(\theta; \cdot)$ at perturbed decision points $\theta + \delta u$ to estimate $\nabla \mathcal{L}$.
- ◇ **Markovian Sample:** $Z_t \sim \Pi_{\theta}$ ✗.
- * cannot draw samples directly from $\Pi_{\theta} \rightarrow$ sample reweighting using the forgetting factor λ .



DFO(λ) Algorithm

- ◇ **Idea:** Construct zero-th order $\mathcal{O}(\delta)$ -biased gradient estimator for $\mathcal{L}(\theta)$ as g_{δ} , to avoid evaluating a priori **unknown** $\Pi_{\theta}(\cdot)$
- $$g_{\delta}(\theta; u, Z) := \frac{d}{\delta} \ell(\check{\theta}; Z) u, \text{ with } \check{\theta} := \theta + \delta u, Z \sim \Pi_{\check{\theta}}(\cdot), u \sim \text{Unif}(\mathbb{S}^{d-1})$$
- ◇ unbiased estimator for $\nabla \mathcal{L}_{\delta}(\theta)$, while $\mathcal{L}_{\delta}(\theta)$ is a smooth approx of $\mathcal{L}(\theta)$.

Two-timescale Derivative Free (DFO(λ)) Algorithm

Outer Loop ($k : 0 \rightarrow T - 1$): Set stepsize δ_k and η_k , inner loop range τ_k
Inner loop ($m : 1 \rightarrow \tau_k$):
 Deploy $\check{\theta}_k^{(m)} = \theta_k^{(m)} + \delta_k u_k$, Sample $Z_k^{(m)} \sim \mathbb{T}_{\check{\theta}_k^{(m)}}(Z_k^{(m-1)}, \cdot)$
 Update $\theta_k^{(m+1)} = \theta_k^{(m)} - \eta_k \lambda^{\tau_k - m} g_{\delta_k}(\theta_k^{(m)}, u_k, Z_k^{(m)}) \rightarrow$ **forgetting factor**
End inner loop: $Z_{k+1} \leftarrow Z_k^{(\tau_k)}, \theta_{k+1} \leftarrow \theta_k^{(\tau_k+1)}$.
Output: θ_s , where $s \sim \text{Uniform}(\{0, 1, \dots, T\})$.

- * **Highlights:** (i) two-timescales step sizes, (ii) make use of every sample, (iii) trade-off between sample accumulation and MC mixing time via λ .

Main Results

- ◇ A1. $\|\nabla \mathcal{L}(\theta) - \nabla \mathcal{L}(\theta')\| \leq L \|\theta - \theta'\|$ A2. (Bounded loss) $|\ell(\theta; z)| \leq G$
- ◇ A3a. (Lipschitz distribution map) $\delta_{TV}(\theta, \theta') \leq L' \|\theta - \theta'\|$
- ◇ A3b. (L_1 -sensitivity) $|\ell(\theta, z) - \ell(\theta, z')| \leq L_0 \|z - z'\|$, $W_1(\Pi_{\theta}, \Pi_{\theta'}) \leq L_1 \|\theta - \theta'\|$
- ◇ A4. (Geometric mixing) $\delta_{TV}(\mathbb{P}_{\theta}(Z_k \in \cdot | Z_0 = z), \Pi_{\theta}) \leq M \rho^k$.
- ◇ A5. (Smooth Markov kernel) $\delta_{TV}(\mathbb{T}_{\theta}(z, \cdot), \mathbb{T}_{\theta'}(z, \cdot)) \leq L_2 \|\theta - \theta'\|$

Theorem 1: Using step sizes $\eta_k \propto k^{-2/3}$, $\delta_k \propto k^{-1/6}$, $\tau_k \propto \log k$, the iterate of DFO(λ) satisfies $\frac{1}{1+T} \sum_{k=0}^T \mathbb{E} \|\nabla \mathcal{L}(\theta_k)\|^2 \leq \mathcal{O}(d^{2/3}/T^{1/3})$

- ◇ **ϵ -stationary:** above metric achieves ϵ -target acc. after $\mathcal{O}(d^2/\epsilon^3)$ iter.
- ◇ **Sample complexity:** $S_{\epsilon} = \mathcal{O}(d^2/\epsilon^3) \leftarrow$ worse than $\mathcal{O}(d/\epsilon^2)$.
- ◇ **Estimator (I):** prior two point estimator $g_{2\text{pt-I}}$ [Ghadimi & Lam, 2013] $g_{2\text{pt-I}} := \frac{d}{\delta} [\ell(\theta + \delta u; Z) - \ell(\theta; Z)] u \rightarrow$ **biased ✗** since $Z \sim \Pi_{\theta + \delta u}$, which is unique feature of decision-dependent sample distribution.
- ◇ **Estimator (II):** $g_{2\text{pt-II}} := \frac{d}{\delta} [\ell(\theta + \delta u; Z_1) - \ell(\theta; Z_2)] u \rightarrow$ **unbiased** Same variance $\mathbb{E} \|g_{2\text{pt-II}}\|^2 = \Omega(1/\delta^2)$, but **higher** sampling overhead ✗

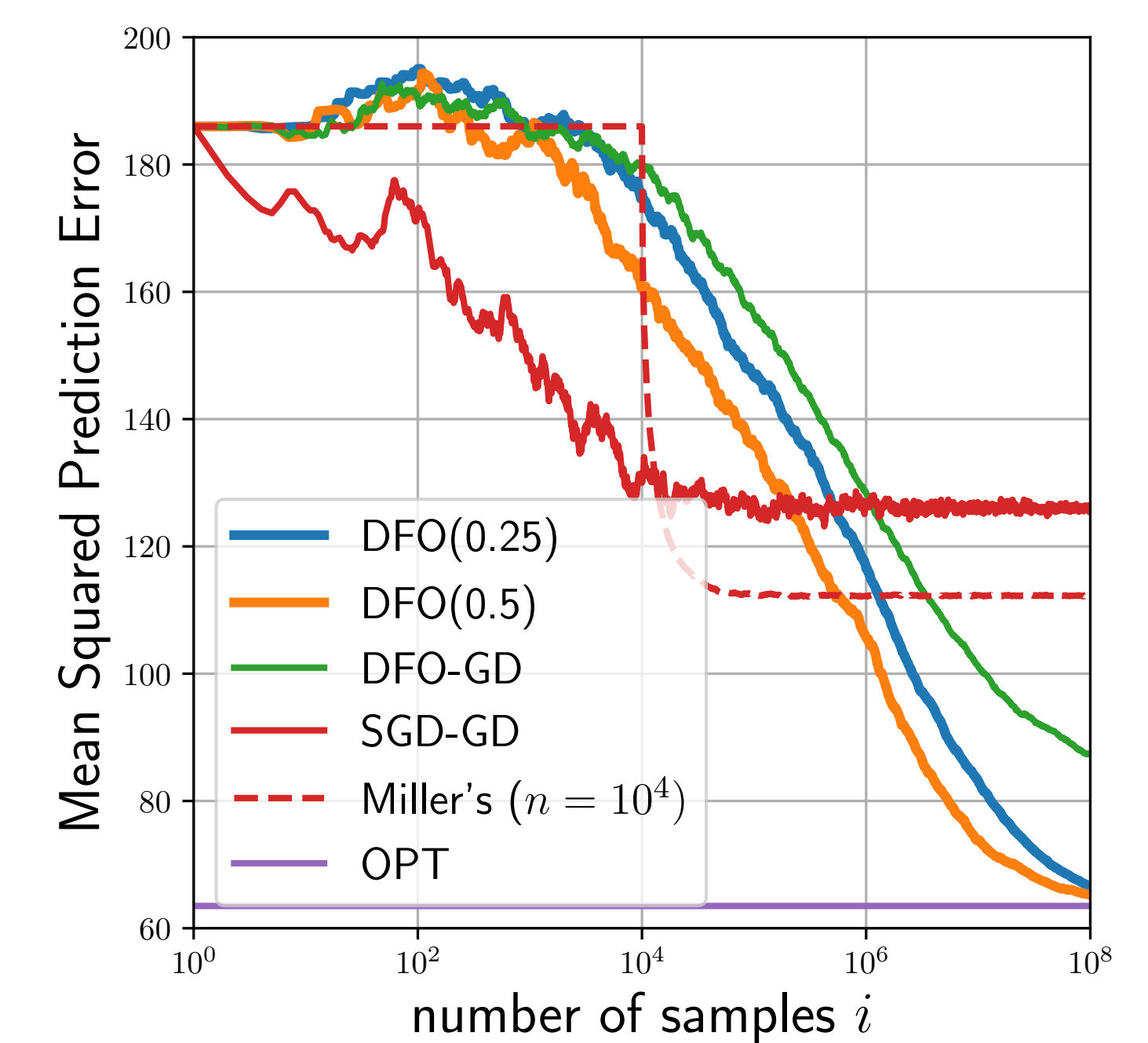
Numerical Experiments

Markovian Regression —

- ◇ **Quadratic Loss:** $\ell(\theta; x, y) = (\langle x, \theta \rangle - y)^2$
- ◇ **AR Model:** $(\tilde{X}_t, \tilde{Y}_t) = (1 - \gamma)(\tilde{X}_{t-1}, \tilde{Y}_{t-1}) + \gamma(X_t, Y_t)$, where γ controls the mixing rate of the Markov chain.
- ◇ **Stationary samples:** (X_t, Y_t) is drawn according to
$$\begin{cases} X_t \sim \mathcal{N}(0, \frac{2-\gamma}{\gamma} \sigma_1^2 I), \\ Y_t | X_t \sim \mathcal{N}(\langle x_t + \kappa \theta_{t-1}, \theta_0 \rangle, \frac{2-\gamma}{\gamma} \sigma_2^2), \end{cases}$$
 where $\gamma = 0.25$, $\kappa = 1/\|\theta_0\|$, $\sigma_1 = \sigma_2 = 1$.
- ◇ **Goal:** Comparison of 4 state-of-the-art algorithms:
 - ◇ SGD with greedy deployment from [Mendler et al., 2020],
 - ◇ Two-Phase algorithm from [Miller et al., 2021]:
 - ◇ (Phase I) Estimate distribution map Π_{θ}
 - ◇ (Phase II) Minimize finite-sample approx. of $\mathcal{L}(\theta)$,
 - ◇ DFO-GD (no burn-in phase),
 - ◇ Our DFO(λ) with $\lambda \in \{0.25, 0.5\}$.

Observations

- ◇ ✗ DFO/SGD-GD fail to find a stationary solution to $\mathcal{L}(\theta)$.
- ◇ ✗ Two-Phase algorithm fail neither even with 10^4 (Markovian) samples gathered in the first phase.
- ◇ ✓ Compared to above algorithms, DFO(λ) converges to a near-optimal solution after reasonable samples.



Reference

- ◇ Perdomo, Juan, et al. *Performative prediction*, ICML 2020.
- ◇ Mendler-Dünner, et al. *Stochastic optimization for performative prediction* NeurIPS 2020.
- ◇ Miller, et al. *Outside the echochamber: Optimizing the performative risk*. In ICML 2021.