

Overview/Motivation:

Our work addresses *time-varying stochastic optimization (TV-SO)*, finds *non-parametric* optimality conditions via Malliavin calculus, and devises a *scalable* neural algorithm, competitive with the baselines from *both* complexity and performance perspectives. Non-parametric framework enable us to provide a learning mechanism *insensitive to the parameterization dimension*, the challenge which most of the baselines such as adjoint sensitivity models and path-wise differentiation methods (PDMs) struggle with.

Key Feature

A scalable algorithm *parallel* to PDMs and adjoint sensitivity models for stochastic optimization under distribution drift.

Application

Our work has applications in machine learning settings where parameters / variables undergo *distribution drift* as well as in *non-stationary environments*. Evaluations on *least-square recovery* and *logistic regression* problems show the effectiveness of our method against baselines, achieving a *10-fold* smaller optimality distance than the stochastic gradient based baseline, and having *extremely less runtime* than the conventional neural-based baseline (PDM).

Problem Statement

Solve the stochastic counterpart of the time-varying optimization (TV-O) problem: (TV-O) : $\min_{\mathbf{x}_t \in [0, T]} f(\mathbf{x}_t, t), \forall t \in [0, T]$. Specifically, solving time-varying stochastic optimization (TV-SO) problem: $\min_{\mathbf{X} = (\mathbf{x}_t \in \mathbb{R}^n)_{t \in [0, T]}} \mathbb{E}_{(\mathbf{z}, \mathbf{X}) \sim \mathbb{P} \otimes \mathbb{Q}} \{f(\mathbf{x}_t, \xi_t, t)\}, \forall t \in [0, T]$, where (i) $\mathbf{z} = (\xi_t \in \mathbb{R}^m)_{t \in [0, T]}$ Malliavin differentiable *known process* under law \mathbb{P} , (ii) $\mathbf{X} = (\mathbf{x}_t \in \mathbb{R}^n)_{t \in [0, T]}$ Malliavin differentiable *decision process* under law \mathbb{Q} , (iii) $f(\cdot, \cdot, \cdot): \mathbb{R}^m \times \mathbb{R}^n \times [0, T] \rightarrow \mathbb{R}$ a square-integrable function representing a time-varying measure.

Theory

Theorem: Consider problem TV-SO under Assumptions (i)–(iii), with $\mathbf{X} = (\mathbf{x}_t)_{t \in [0, T]}$ and $\mathbf{z} = (\xi_t)_{t \in [0, T]}$ being \mathcal{F} -adapted, Malliavin-differentiable, and square-integrable processes. Then TV-SO is optimized at a given instant t if the following necessary conditions hold:

$$\mathbb{E}\{\nabla_{\mathbf{x}}^2 f(\mathbf{x}_t, \xi_t, t) D_s \mathbf{x}_t + \nabla_{\mathbf{x} \xi} f(\mathbf{x}_t, \xi_t, t) D_s \xi_t\} = \mathbf{0}, \quad \forall s \leq t$$

$$\mathbb{E}\{\nabla_{\mathbf{x}} f(\mathbf{x}_t, \xi_t, t)\} = \mathbf{0}$$

where $D_s \mathbf{x}_t$ (respectively, $D_s \xi_t$) is the Malliavin derivative of \mathbf{x}_t (respectively, ξ_t) with respect to Wiener process \mathbf{w}_t .

Proposition: Consider the Malliavin derivative $D_s \mathbf{x}_t$ for the SDE $d\mathbf{x}_t = \mathbf{a}(\mathbf{x}_t, t)dt + \sum_{l=1}^d \mathbf{b}_l(\mathbf{x}_t, t)dw_t^l$. Then, $D_s \mathbf{x}_t = \Gamma_{s,t} \mathbf{B}(\mathbf{x}_s, s)$ for $s \leq t$, where $\mathbf{B}(\mathbf{x}_s, s) = [\mathbf{b}_1, \dots, \mathbf{b}_d](\mathbf{x}_s, s)$ and s -to- t stochastic flow $\Gamma_{s,t} \in \mathbb{R}^{n \times n}$ represents a linear SDE, adapted to the filtration \mathcal{F}_s , given by

$$d\Gamma_{s,t} = \left(\sum_{l=1}^d J_x \mathbf{b}_l(\mathbf{x}_t, t) dw_t^l + J_x \mathbf{a}(\mathbf{x}_t, t) dt \right) \Gamma_{s,t}, \quad \Gamma_{s,s} = \mathbf{I}, \quad J_x[\cdot] \rightarrow \text{Jacobian operator}$$

Comparison of Optimality Conditions

TV-O (literature): $\nabla_{\mathbf{x}}^2 f \frac{d\mathbf{x}_t}{dt} + \frac{\partial}{\partial t} \nabla_{\mathbf{x}} f = \mathbf{0}$

TV-SO (our problem): $\mathbb{E}\{\nabla_{\mathbf{x}}^2 f D_s \mathbf{x}_t + \nabla_{\xi} \nabla_{\mathbf{x}} f D_s \xi_t\} = \mathbf{0}$

A Stochastic Path Follower (SPF) algorithm to solve TV-SO

(i) **Representation:** represent decision process \mathbf{x}_t using a neural SDE

$$d\mathbf{x}_t^\theta = \mathbf{a}(\mathbf{x}_t^\theta, t; \theta)dt + \sum_{l=1}^d \mathbf{b}_l(\mathbf{x}_t^\theta, t; \theta)dw_t^l$$

(ii) **Forward Pass⁽¹⁾:** simulate the process \mathbf{x}_t^θ using any SDE-solver, given NN parameters θ

(iii) **Forward Pass⁽²⁾:** simulate the Malliavin derivatives, $(D_s \mathbf{x}_t^\theta, D_s \xi_t)$, using any SDE-solver and based on the s -to- t stochastic flow $\Gamma_{s,t}$

(iv) **Loss Evaluation:** include both the optimality conditions by the following energy-functional criterion as the learning loss:

$$\mathcal{L}(\theta; \mathbf{x}, D\mathbf{x}) = \int_0^T \left(\|\mathbb{E} \nabla_{\mathbf{x}} f\|^2 + \int_0^t \|\mathbb{E}\{\nabla_{\mathbf{x}}^2 f D_s \mathbf{x}_t^\theta + \nabla_{\mathbf{x} \xi} f D_s \xi_t\}\|^2 ds \right) dt$$

(v) **Update:** update the parameters of the neural SDE using a convergence-guaranteed Adam-type algorithm.

Algorithm 1: Stochastic Path Follower (SPF)

Goal: Solving stochastic optimization problem TV-SO .

Input: Neural drift and diffusion functions $\mathbf{a}(\cdot, t, \theta): \mathbb{R}^{n \times p} \times [0, T] \times \Theta \rightarrow \mathbb{R}^{n \times p}$, $\mathbf{B}(\cdot, t, \theta): \mathbb{R}^{n \times p} \times [0, T] \times \Theta \rightarrow \mathbb{R}^{n \times p \times d}$, initial process $\mathbf{x}_0 \in \mathbb{R}^{n \times p}$, number of sample paths P , number of episodes N_{eps} .

Output: Diffusion process \mathbf{X} with optimal solution θ of TV-SO.

for $E = 1$ to N_{eps} do

 // Simulate the neural SDE to obtain sample paths:

$(\mathbf{x}_t, \mathbf{w}_t) \leftarrow \text{SDEsolve}(\mathbf{a}, \mathbf{B}, \mathbf{x}_0, \theta)$;

 // Compute Malliavin derivative $D_s \mathbf{x}_t$ by simulating the SDE

 // of the s -to- t stochastic flow $\Gamma_{s,t}$

$D_s \mathbf{x}_t \leftarrow \text{SDEsolve}(J_x \mathbf{a}, J_x \mathbf{B}, \Gamma_{s,s})$;

 // Evaluate the energy-functional loss $\mathcal{L}(\theta; \mathbf{x}, D\mathbf{x})$

$\mathcal{L}(\theta) \leftarrow \mathcal{L}(\theta; \mathbf{x}, D\mathbf{x})$;

 // Update parameters of neural SDE:

$\theta \leftarrow \text{AdamOptimizer}(\theta, \mathcal{L}(\theta))$;

end

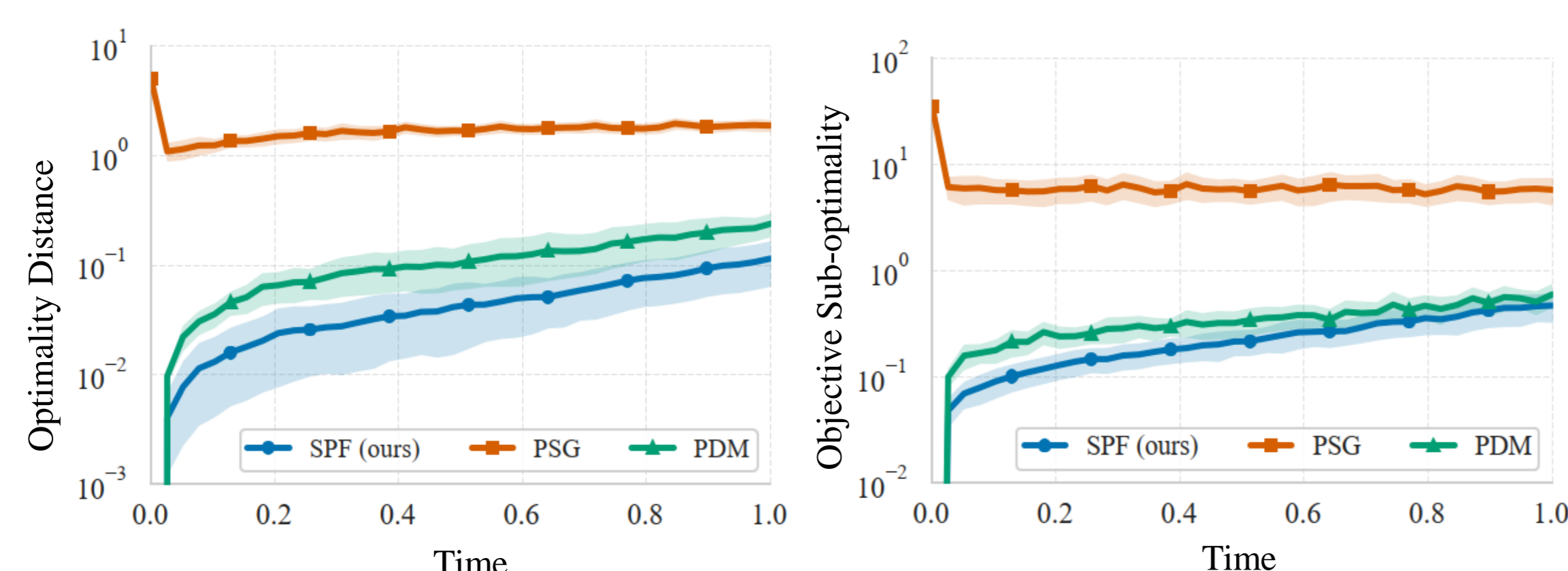
Experiments

Approaches	Stochastic Path Follower SPF (ours)	Scalable neural algorithm; needs $n + n^2$ SDEs for solution (notably independent of parametrization)
	Path-wise Differentiation Method (PDM) [Tze19]	Neural baseline; needs $n + n\phi$ SDEs for solution with ϕ parameterization dimension
	Proximal Stochastic Gradient (PSG) [Cut23]	Neural-free gradient-based approach

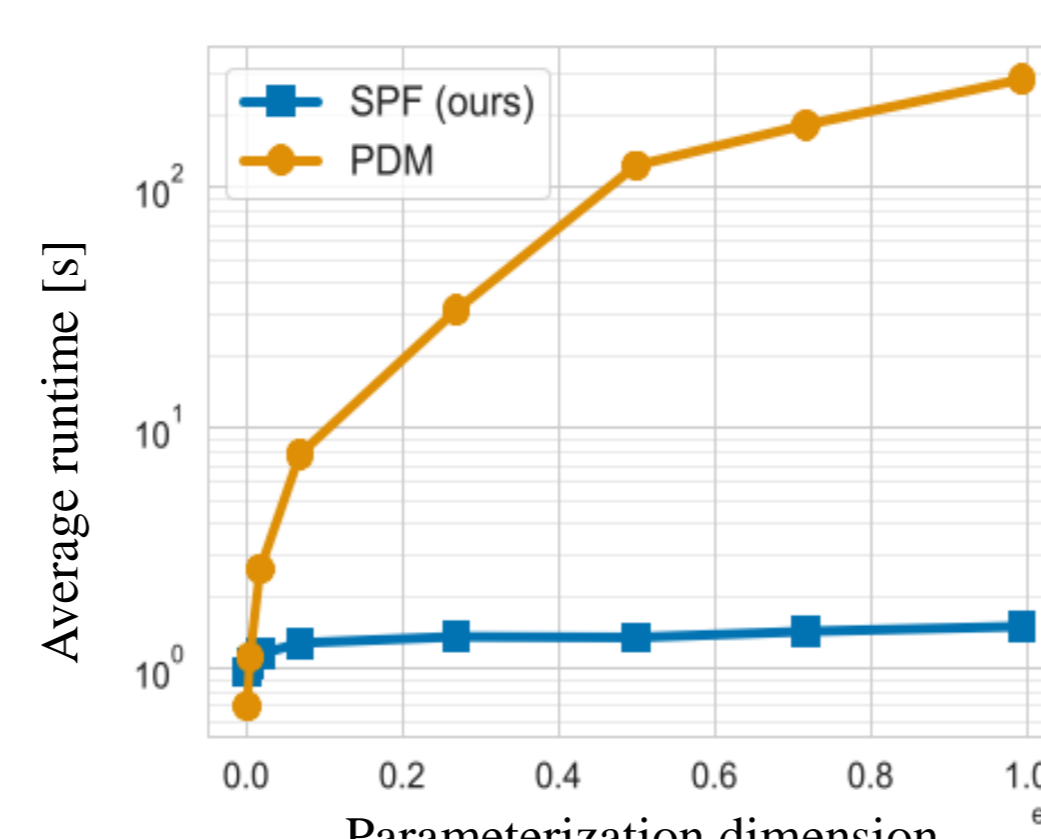
Least-squares recovery with distribution drift:

$$\min_{\mathbf{X}} \mathbb{E} \frac{1}{2} \|\mathbf{A} \mathbf{x}_t - \mathbf{n}_t\|^2, \quad \mathbf{A} \in \mathbb{R}^{n \times d}, \mathbf{n}_t \in \mathbb{R}^n, \mathbf{x}_t \in \mathbb{R}^d$$

\mathbf{A} : measurement matrix, $\mathbf{n}_t \sim \mathcal{N}(\mathbf{A} \mathbf{x}_t^*, \mathbf{C})$: observed parameter, $\mathbf{C} = \sigma^2 \mathbf{I}$
 \mathbf{x}_t^* follows an Ornstein–Uhlenbeck (OU) process



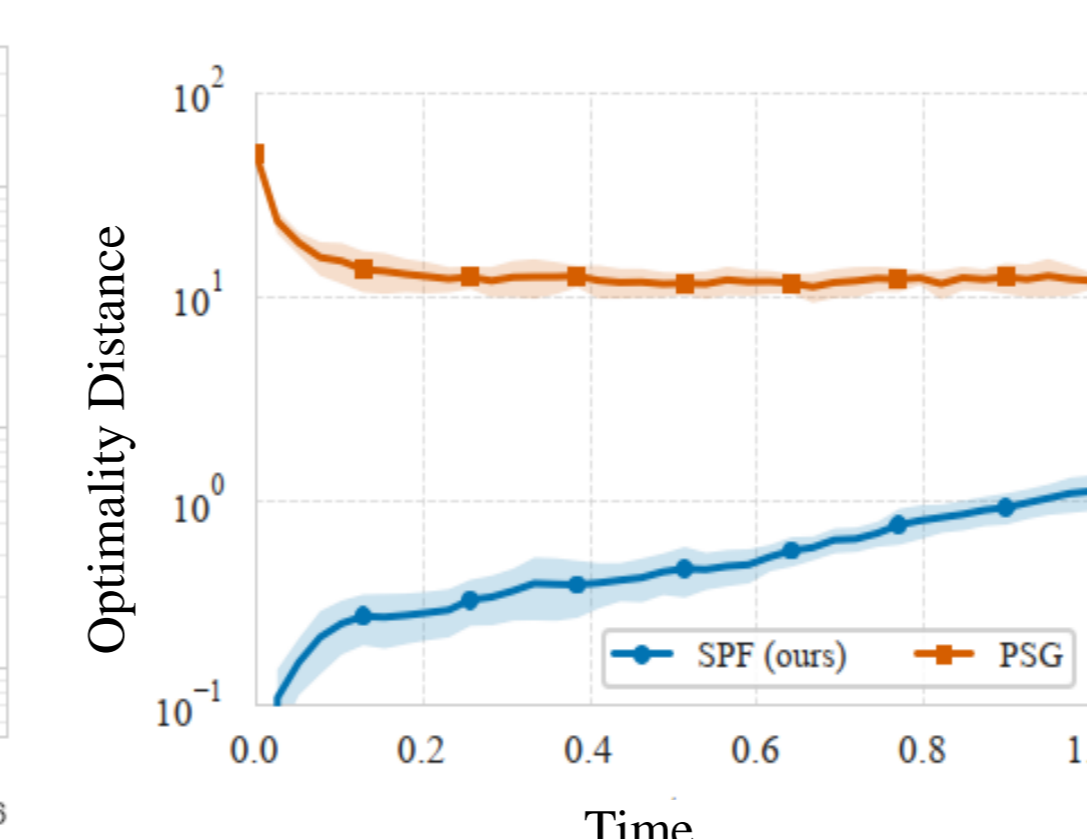
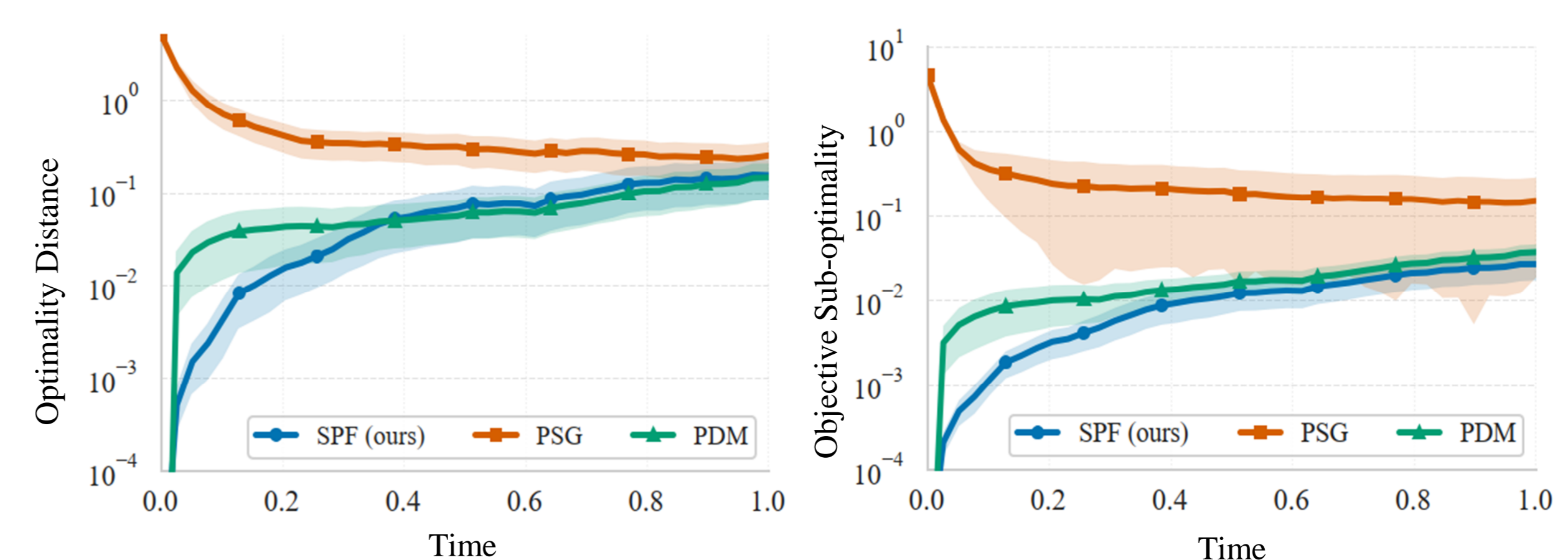
High-dimensional evaluation on least-squares recovery ($n = 100, d = 50$)



Logistic regression with distribution drift:

$$\min_{\mathbf{X}} \mathbb{E} \frac{1}{n} \left(\sum_{i=1}^n \log(1 + \exp(\mathbf{a}_i^T \mathbf{x}_t)) - \mathbf{b}_t^T \mathbf{A} \mathbf{x}_t \right), \quad \mathbf{A} = \{\mathbf{a}_i\}_{i=1}^n \in \mathbb{R}^{n \times d}, \mathbf{b}_t \in \mathbb{R}^n, \mathbf{x}_t \in \mathbb{R}^d$$

\mathbf{A} : measurement matrix, $\mathbf{b}_t = 1/(1 + e^{-\mathbf{y}_t})$: soft labels with \mathbf{y}_t following an OU process



Quick Access

