

Towards Robust Learning to Optimize with Theoretical Guarantee

Qingyu Song[†], Wei Lin[†], Juncheng Wang[‡], Hong Xu[†]

[†]The Chinese University of Hong Kong, [‡]Hong Kong Baptist University



香港中文大學
The Chinese University of Hong Kong



 香港浸會大學
HONG KONG BAPTIST UNIVERSITY

What is learning to optimize (L2O)?

- Optimization Problem

$$\min_x f(x), \\ x \in R^n, f: R^n \rightarrow R$$

Human-Designed Algorithm

Gradient
Descent



Learning to Optimize

- Benefits

- Better optimality (potential).
- Better convergence/efficiency [1].

ML/DL
Model

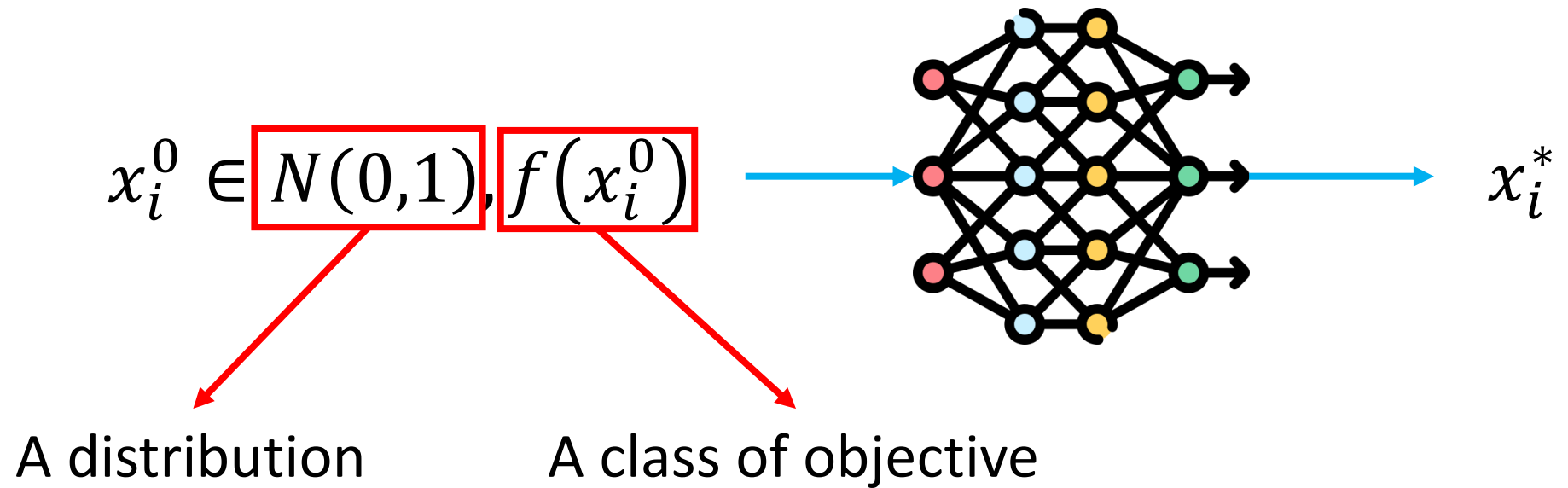
+

Training

How does L2O work?

- Workflow of L2O (Inference)

Given initial point x_i^0



L2O's Failure in Out-of-Distribution (OOD) Scenarios

① x OOD

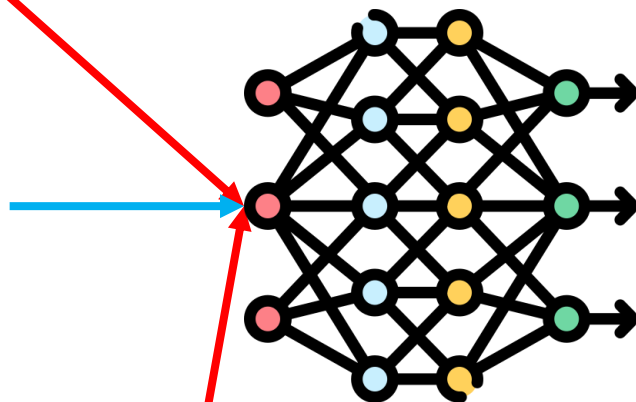
$$x^0 \in N(\mathbf{10}, 1), f(x)$$

In-Distribution (InD)

$$x^0 \in N(\mathbf{0}, 1), f(x)$$

② f OOD

$$x^0 \in N(\mathbf{0}, 1), f = f(x - \mathbf{10})$$



$$x^* \in D'$$

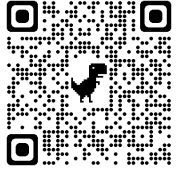
$$\star x^* \neq x^*$$

$$x^* \in D$$

$$\star x^* \neq x^* + 10$$

$$x^* \in D''$$

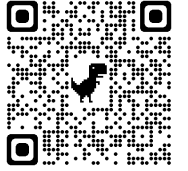
Theoretical Convergence Analysis



- Convergence of Single-Iteration (Smooth Case)

$$\begin{aligned} & F'(x_k + s_k) - F'(x_{k-1} + s_{k-1}) \\ \leq & \boxed{-\frac{\|\nabla f'(x_{k-1} + s_{k-1})\|^2}{2L}} \quad \text{Convergence of Gradient-Descent} \\ & + L \|\text{diag}(\mathbf{J}_{1,k-1} s') \nabla f'(x_{k-1} + s_{k-1})\|^2 \\ & + L \left\| \frac{\nabla f'(x_{k-1} + s_{k-1}) - \nabla f(x_{k-1})}{2L} - \mathbf{J}_{2,k-1} s' \right\|^2. \end{aligned} \quad \left. \vphantom{\frac{\nabla f'(x_{k-1} + s_{k-1}) - \nabla f(x_{k-1})}{2L}} \right\} \text{Deterioration w.r.t. OOD}$$

Theoretical Convergence Rate Analysis



- Convergence Rate (Smooth Case)

$$\begin{aligned} & \min_{k=1, \dots, K} F'(x_k + s_k) - F'(x^* + s^*) \\ & \leq \frac{L}{2} \|x_0 - x^* + s_0 - s^*\|^2 - \frac{L}{2} \|x_K - x^* + s_K - s^*\|^2 \\ & \quad + \frac{L}{K} \sum_{k=1}^K (x_k + s_k - x^* - s^*)^\top \\ & \quad \left(x_k + s_k - \left(x_{k-1} + s_{k-1} - \frac{\nabla f'(x_{k-1} + s_{k-1})}{L} \right) \right). \end{aligned}$$

Convergence Rate of
Gradient-Descent

Deterioration w.r.t. OOD

Convergence Improvement



- Upper Bound Relaxation

$$\begin{aligned} & F'(x_k + s_k) - F'(x_{k-1} + s_{k-1}) \\ & \leq - \frac{\|\nabla f'(x_{k-1} + s_{k-1})\|^2}{2L} \\ & \quad + \frac{\|\nabla f'(x_{k-1} + s_{k-1}) - \nabla f(x)\|^2}{2L} \\ & \quad + (LC_1^2 n \|\nabla f'(x_{k-1} + s_{k-1})\|^2 + 2LC_2^2 n) \boxed{\|s'\|^2}. \end{aligned}$$

OOD vector,
NN's input feature

- Improve upper bound: Magnitude reduction.
 - Our approach: Input Feature Simplification.

A New L2O Model with Gradient-Only Input

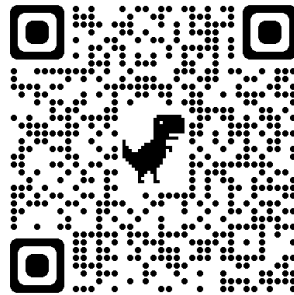
- New Model Formulation Based on [1]

$$x_k = x_{k-1} - \mathbf{R}_k \nabla f(x_{k-1}) - \mathbf{R}_k g_k - \mathbf{Q}_k v_{k-1} - b_{1,k},$$

$$v_k = (\mathbf{I} - \mathbf{B}_k) G_k + \mathbf{B}_k G_{k-1} - b_{2,k},$$

$$G_k := \mathbf{R}_k^{-1} (x_{k-1} - x_k - \mathbf{Q}_k v_{k-1} - b_{1,k}),$$

- Learn \mathbf{R} , \mathbf{Q} , \mathbf{B} . Details at



Empirical Outperformance

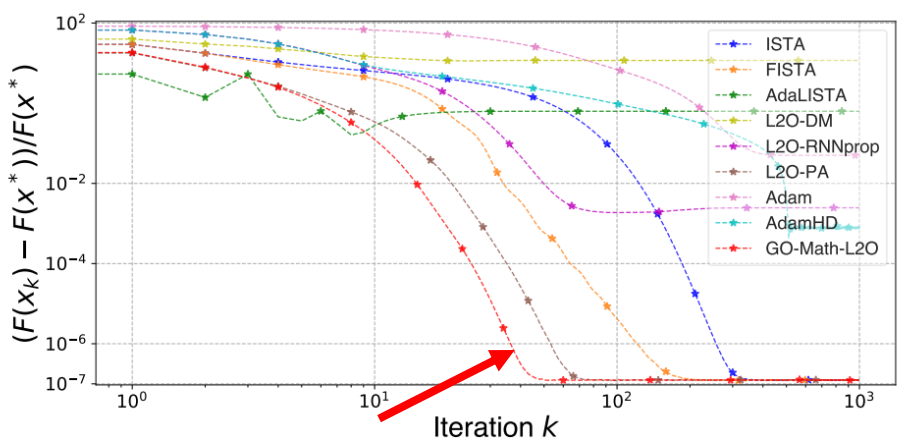


Figure 1. LASSO Regression: InD.

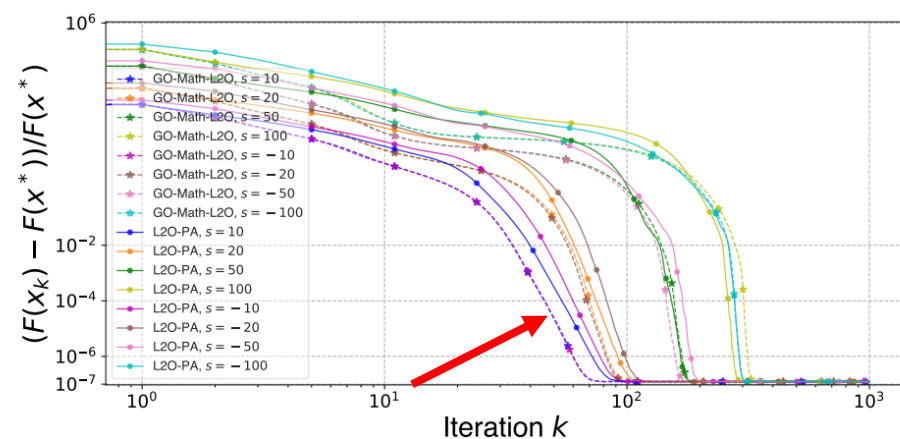


Figure 3. LASSO Regression: OOD by Trigger 1.

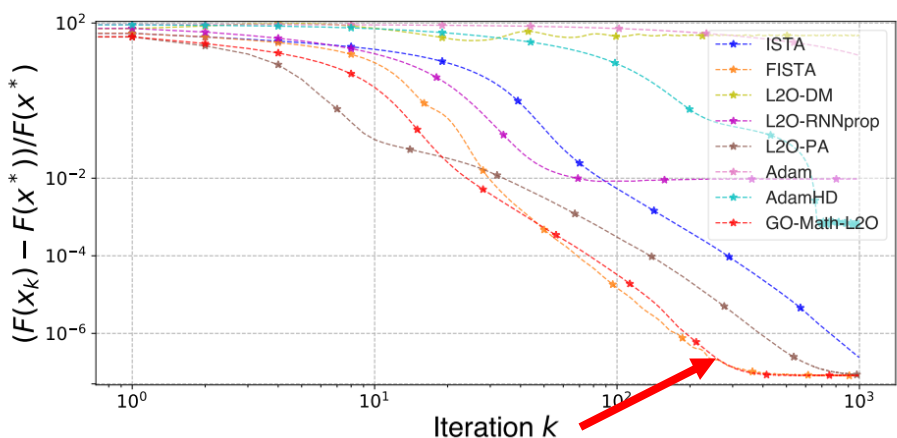


Figure 2. LASSO Regression: Real-World OOD.

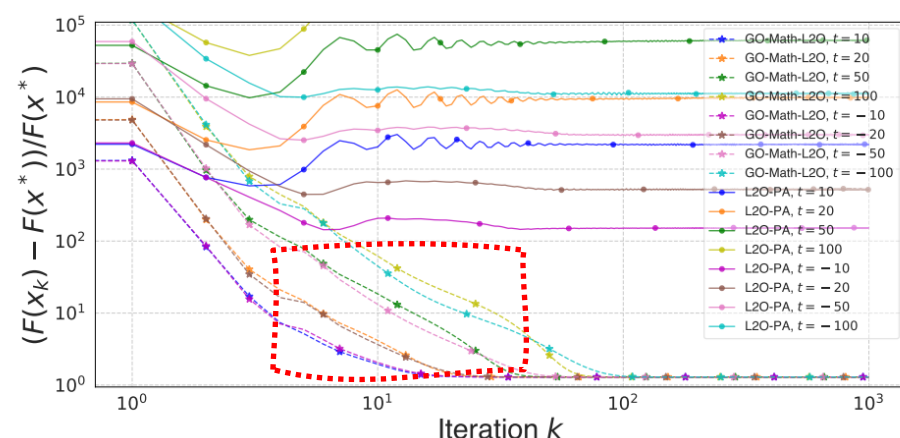


Figure 4. LASSO Regression: OOD by Trigger 2.

Empirical Outperformance

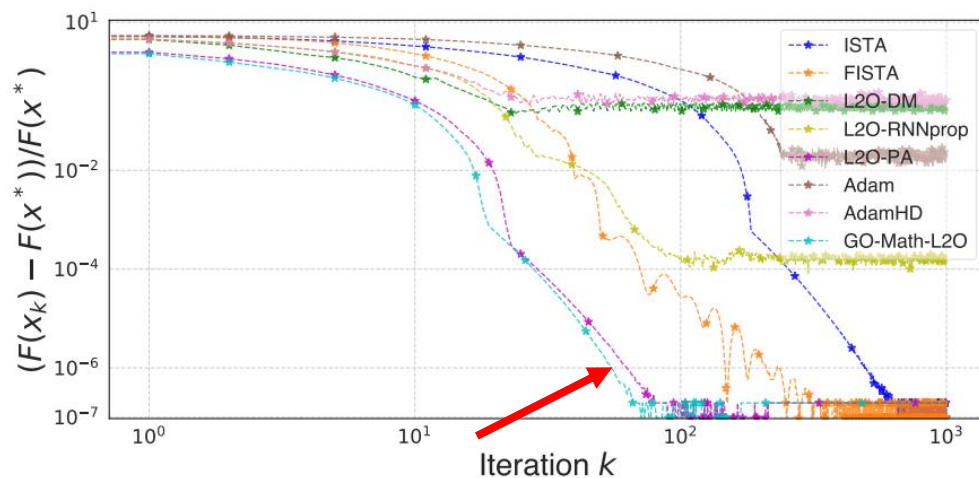


Figure 14. Logistic Regression: Real-World Ionosphere Dataset.

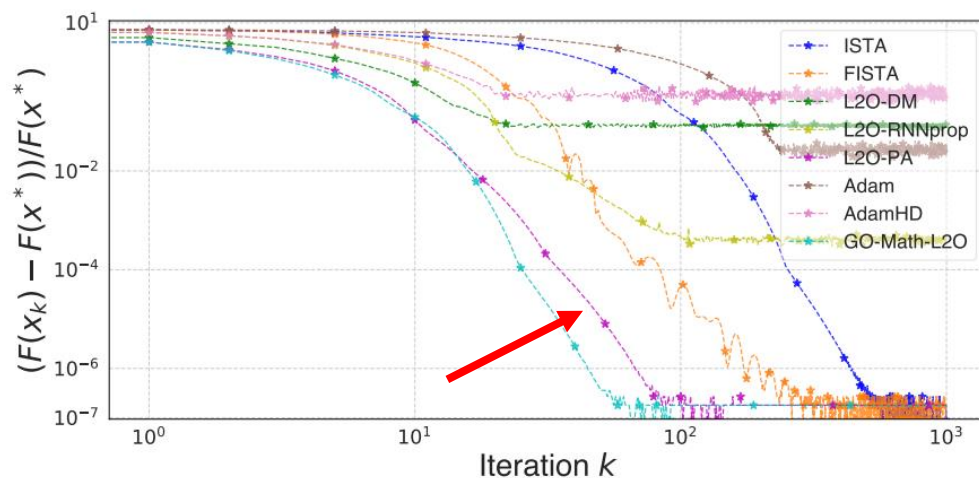


Figure 15. Logistic Regression: Real-World Spambase Dataset.

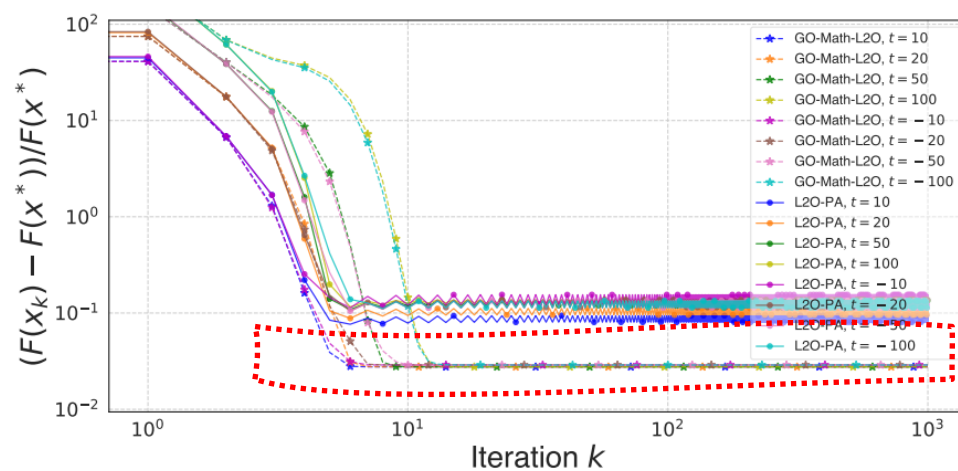


Figure 17. Logistic Regression: OOD by Trigger 2.

Thank You!

Project at [NetX-lab/GoMathL2O-Official \(github.com\)](https://github.com/NetX-lab/GoMathL2O-Official)

Paper



Code



香港中文大學
The Chinese University of Hong Kong



香港浸會大學
HONG KONG BAPTIST UNIVERSITY