



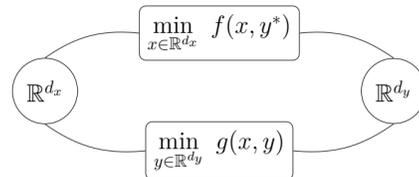
Bilevel optimization problem (BiO)

General form

$$\begin{aligned} \min_{x \in \mathbb{R}^{d_x}} \quad & \varphi(x) := f(x, y^*(x)) \\ \text{s. t.} \quad & y^*(x) = \arg \min_{y \in \mathbb{R}^{d_y}} g(x, y) \end{aligned}$$

► $f, g: \mathbb{R}^{d_x} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}$, differentiable, g is **strongly convex**

The nested structure **ouples** the upper level and lower level



Approximate implicit differentiation (AID)

Hyper-gradient computation

$$\begin{aligned} \nabla \varphi(x) &= \nabla_x f(x, y^*(x)) + \nabla_x y^*(x)^\top \nabla_y f(x, y^*(x)) \\ &= \nabla_x f(x, y^*) - \nabla_{xy}^2 g(x, y^*) [\nabla_{yy}^2 g(x, y^*)]^{-1} \nabla_y f(x, y^*) \end{aligned}$$

Main difficulties

► Solving the lower-level problem to obtain $y^*(x)$

★ Approximating the **Hessian inverse vector product**

$$v^*(x) := [\nabla_{yy}^2 g(x, y^*(x))]^{-1} \nabla_y f(x, y^*(x))$$

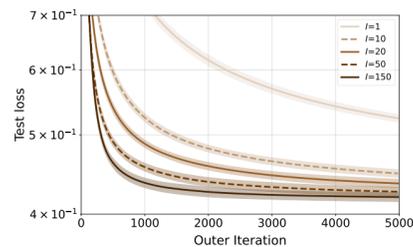
Vanilla update rule

$$\mathbf{1} \times : x^+ = x - \beta (\nabla_x f(x, y) - \nabla_{xy}^2 g(x, y) [\nabla_{yy}^2 g(x, y)]^{-1} \nabla_y f(x, y))$$

$$\mathbf{N} \times : y^+ = y - \alpha \nabla_y g(x, y)$$

Motivation

An observation



The more **accurate** v^* ,
the more **enhanced** descent!

How to tackle v^* in BiO

► **Approximation** principle

► **Amortization** principle

How to adhere to these two principles

► Subspace methods for **efficient Approximation** [Yuan, 2014]

► Subspace iteration for **reasonable Amortization** [Yuan, 1995]

A subspace perspective

Geometric interpretation

$$\min_{v \in \mathcal{S}_k} m_k(v) := \frac{1}{2} v^\top \nabla_{yy}^2 g(x_k, y_k) v - \nabla_y f(x_k, y_k)^\top v$$

$$\mathcal{S}_k = \mathbb{R}^{d_y} \implies v_k = \nabla_{yy}^2 g(x_k, y_k)^{-1} \nabla_y f(x_k, y_k) := A_k^{-1} b_k$$

Low-dimensional subspace?

Krylov subspace [Krylov, 1931] provides a good estimate for $A^{-1}b$ [Carmon and Duchi, 2018]

$$\mathcal{K}_N(A, b) := \text{span} \{b, Ab, A^2b, \dots, A^{N-1}b\}$$

Constructing 2-D subspace: SubBiO

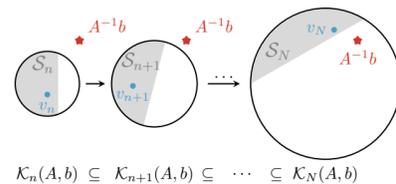
Given $v_n \in \mathcal{K}_n(A, b)$, an initial approximation of $A^{-1}b$

$$v_n = \sum_{i=0}^{n-1} c_i (I - \eta A)^i b \approx A^{-1}b$$

Recursively, we can choose

$$v_{n+1} \in \mathcal{S}_{n+1} := \text{span} \{b, (I - \eta A)v_n\} \subseteq \mathcal{K}_{n+1}(A, b)$$

SubBiO



$$\mathcal{K}_n(A, b) \subseteq \mathcal{K}_{n+1}(A, b) \subseteq \dots \subseteq \mathcal{K}_N(A, b)$$

$$\mathbf{1} \times : \mathcal{S}_k := \text{span} \{b_k, (I - \eta A_k)v_{k-1}\}$$

$$\mathbf{1} \times : v_k := \arg \min_{v \in \mathcal{S}_k} m_k(v) := \frac{1}{2} v^\top A_k v - b_k^\top v$$

$$\mathbf{1} \times : x_{k+1} = x_k - \beta (\nabla_x f(x_k, y_k) - \nabla_{xy}^2 g(x_k, y_k) v_k)$$

$$\mathbf{1} \times : y_{k+1} = y - \alpha \nabla_y g(x_k, y_k)$$

Computational complexity

► $\min_{z \in \mathbb{R}^2} \frac{1}{2} z^\top (S_k^\top A_k S_k) z - b_k^\top S_k z: O(2n^2)$

► $A_k v_{k-1}: O(n^2), \nabla_{xy}^2 g(x_k, y_k) v_k: O(n^2)$

SubBiO costs $O(4n^2)$ per iteration

Dynamic Lanczos process in BiO

Core principles

► Maintain an **orthogonal basis** $Q_j = [q_1, \dots, q_j]$ approximating $\mathcal{K}_j(A_j, b_j)$, $\mathcal{S}_k = \text{span}(Q_j)$

► Keep the (approximate) projection matrix T_j **tridiagonal**

► Dynamically solve quadratic subproblems: $v_k := \arg \min_{v \in \mathcal{S}_k} m_k(v) := \frac{1}{2} v^\top A_k v - b_k^\top v$

Adapt standard Lanczos process in BiO

$$\begin{aligned} u_j &= A_j q_j - \beta_j q_{j-1}, & \alpha_j &= q_j^\top u_j & \omega_j &= u_j - \alpha_j q_j \\ \beta_{j+1} &= \|\omega_j\| & q_{j+1} &= \omega_j / \beta_{j+1} \end{aligned} \quad T_j = \begin{pmatrix} & & & & 0 \\ & & & & \beta_j \\ & & T_{j-1} & & \\ & & & & \beta_j \\ 0 & & & & \alpha_j \end{pmatrix}$$

Lanczos process is inherently unstable [Paige, 1980; Meurant and Strakoš, 2006]

Two “Res” strategies: LancBiO

Restart mechanism

► Restart subspaces each m steps

► Mitigate the accumulation of difference among $\{A_1, \dots, A_k\}$

Residual minimization

► Minimal residual subproblems

$$\min_{\Delta v \in \mathcal{S}_k} \|(b_k - A_k \bar{v}) - A_k \Delta v\|^2$$

► Correct current \bar{v} , $v_k = \bar{v} + \Delta v_k$

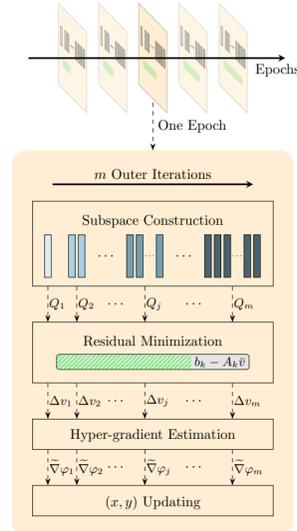
► Collect **historical information**

Computational complexity

⊙ Low-dimensional subproblems

⊙ No cost of Hessian projection

⊙ LancBiO needs $O(2n^2)$ instead of $O(4n^2)$



Theoretical analysis

Proposition

The dynamic Lanczos process in LancBiO with normalized q_1 and α_j, β_j, q_j satisfies

$$A_j^* Q_j = Q_j T_j + \beta_{j+1} q_{j+1} e_j^\top + \delta Q_j, \quad \text{for } j = 1, 2, \dots, m, \quad \text{with } \|\delta q_j\| \leq L_{g_{yy}} \varepsilon_j$$

Theorem

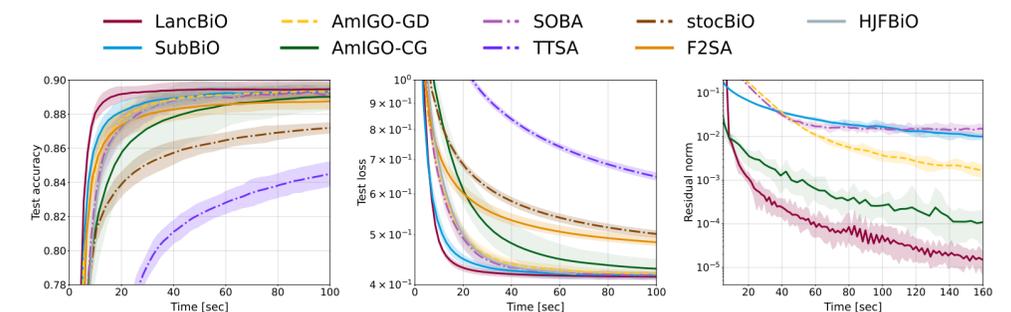
With appropriate step sizes θ for y and λ for x , the iterates $\{x_k\}$ generated by LancBiO satisfy

$$\frac{m}{K(m-m_0)} \sum_{\substack{k=0 \\ (k \bmod m) > m_0}}^K \|\nabla \varphi(x_k)\|^2 = \mathcal{O}\left(\frac{m\lambda^{-1}}{K(m-m_0)}\right)$$

where m is the subspace dimension and $m_0 \sim \Omega(\log m)$

Numerical Experiment

Test on MNIST (pollution rate=0.8)



More experiments are referred to our paper...