

Robust Multi-Objective Bilevel Optimization

Applications in Machine Learning

Pari Ram[†], Alex Gu[¶], Songtao Lu[†], Lily Weng^{†§}

October 18, 2022



$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D} \quad & F(\mathbf{x}) \triangleq \mathbb{E}_{\xi} [f(\mathbf{x}, \mathbf{y}^*(\mathbf{x}); \xi)] \\ \text{subject to} \quad & \mathbf{y}^*(\mathbf{x}) \in \underset{\mathbf{y} \in \mathcal{Y} = \mathbb{R}^d}{\text{arg min}} G(\mathbf{x}, \mathbf{y}) \triangleq \mathbb{E}_{\zeta} [g(\mathbf{x}, \mathbf{y}; \zeta)] \end{aligned}$$

$$\begin{aligned} \min_{x \in \mathcal{X} \subseteq \mathbb{R}^D} F(x) &\triangleq \mathbb{E}_{\xi} [f(x, y^*(x); \xi)] \\ \text{subject to } y^*(x) &\in \operatorname{argmin}_{y \in \mathcal{Y} = \mathbb{R}^d} G(x, y) \triangleq \mathbb{E}_{\zeta} [g(x, y; \zeta)] \end{aligned}$$

Stochastic Objectives

- ▶ **Strongly convex** lower-level (LL) objective $G(x, \cdot)$
- ▶ **Weakly convex** upper-level (UL) objective $F(x)$

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^D} F(\mathbf{x}) &\triangleq \mathbb{E}_{\xi} [f(\mathbf{x}, \mathbf{y}^*(\mathbf{x}); \xi)] \\ \text{subject to } \mathbf{y}^*(\mathbf{x}) &\in \underset{\mathbf{y} \in \mathcal{Y} = \mathbb{R}^d}{\text{argmin}} G(\mathbf{x}, \mathbf{y}) \triangleq \mathbb{E}_{\zeta} [g(\mathbf{x}, \mathbf{y}; \zeta)] \end{aligned}$$

Constraints

- ▶ **Unconstrained** LL problem
- ▶ **Constrained** UL problem

We have n pairs of stochastic objectives $F_i, G_i, i \in [n] \triangleq \{1, \dots, n\}$

$$\begin{aligned} \min_{x \in \mathcal{X} \subseteq \mathbb{R}^D} \quad & \max_{i \in [n]} F_i(x) \triangleq \mathbb{E}_{\xi_i} [f_i(x, y_i^*(x); \xi_i)] \\ \text{subject to} \quad & y_i^*(x) \in \underset{y_i \in \mathcal{Y}_i = \mathbb{R}^{d_i}}{\operatorname{argmin}} G_i(x, y_i) \triangleq \mathbb{E}_{\zeta_i} [g_i(x, y_i; \zeta_i)] \\ & \forall i \in [n] \end{aligned}$$

We have n pairs of stochastic objectives $F_i, G_i, i \in [n] \triangleq \{1, \dots, n\}$

$$\begin{aligned} \min_{x \in \mathcal{X} \subseteq \mathbb{R}^D} \max_{i \in [n]} F_i(x) &\triangleq \mathbb{E}_{\xi_i} [f_i(x, y_i^*(x); \xi_i)] \\ \text{subject to } y_i^*(x) &\in \operatorname{argmin}_{y_i \in \mathcal{Y}_i = \mathbb{R}^{d_i}} G_i(x, y_i) \triangleq \mathbb{E}_{\zeta_i} [g_i(x, y_i; \zeta_i)] \\ &\forall i \in [n] \end{aligned}$$

Multiple objectives

- ▶ Each obj pair has their UL and LL stoc oracles $\xi_i, \zeta_i, i \in [n]$
- ▶ UL variable x **shared across all objectives**
- ▶ Each pair has its **own specific** LL variable $y_i, i \in [n]$

We have n pairs of stochastic objectives $F_i, G_i, i \in [n] \triangleq \{1, \dots, n\}$

$$\begin{aligned} \min_{x \in \mathcal{X} \subseteq \mathbb{R}^D} \max_{i \in [n]} F_i(x) &\triangleq \mathbb{E}_{\xi_i} [f_i(x, y_i^*(x); \xi_i)] \\ \text{subject to } y_i^*(x) &\in \operatorname{argmin}_{y_i \in \mathcal{Y}_i = \mathbb{R}^{d_i}} G_i(x, y_i) \triangleq \mathbb{E}_{\zeta_i} [g_i(x, y_i; \zeta_i)] \\ &\forall i \in [n] \end{aligned}$$

Additional Features

- ▶ **Flexible:** LL vars $\{y_i\}_{i \in [n]}$ can have diff domains ($\mathcal{Y}_i \neq \mathcal{Y}_j$)
- ▶ **Robust:** Shared UL var x optimizes worst-case obj F_i

Learn a robust representation useful for many tasks.

	Problem mapping
$i \in [n]$	Tasks
UL var x	Shared representation network Φ_x params
LL var y_i	Per-task model w_{y_i} params
UL obj $\mathbb{E}_{\xi_i} f_i(x, y_i; \xi_i)$	$\mathcal{L}(w_{y_i} \circ \Phi_x; D_i^{\text{val}})$
LL obj $\mathbb{E}_{\zeta_i} g_i(x, y_i; \zeta_i)$	$\mathcal{L}(w_{y_i} \circ \Phi_x; D_i^{\text{train}}) + \rho \cdot \Omega(y_i)$

Learn a robust representation ensuring fairness across different groups.

	Problem mapping
$i \in [n]$	Demographic groups
UL var x	Shared representation network Φ_x params
LL var y_i	Per-task model w_{y_i} params
UL obj $\mathbb{E}_{\xi_i} f_i(x, y_i; \xi_i)$	$\mathcal{L}(w_{y_i} \circ \Phi_x; D_i^{\text{val}})$
LL obj $\mathbb{E}_{\zeta_i} g_i(x, y_i; \zeta_i)$	$\mathcal{L}(w_{y_i} \circ \Phi_x; D_i^{\text{train}}) + \rho \cdot \Omega(y_i)$

Learn a robust set of hyperparameters (HP) that works well for multiple problems.

	Problem mapping
$i \in [n]$	Different HPO tasks
UL var x	Shared HP x
LL var y_i	Per-task model w_{y_i} params for HP x
UL obj $\mathbb{E}_{\xi_i} f_i(x, y_i; \xi_i)$	$\mathcal{L}(w_{y_i}; D_i^{\text{val}})$
LL obj $\mathbb{E}_{\zeta_i} g_i(x, y_i; \zeta_i)$	$\mathcal{L}(w_{y_i}; D_i^{\text{train}}), w_{y_i} \in \mathcal{F}(x)$

Learn a shared SuperNet such that different application specific sub-networks obtained via Differentiable Architecture Search (DARTS) have robust performance.

	Problem mapping
$i \in [n]$	Different applications
UL var x	SuperNet params W_x
LL var y_i	Per-application subnetwork w_{y_i} params
UL obj $\mathbb{E}_{\xi_i} f_i(x, y_i; \xi_i)$	$\mathcal{L}(w_{y_i}; D_i^{\text{val}})$
LL obj $\mathbb{E}_{\zeta_i} g_i(x, y_i; \zeta_i)$	$\mathcal{L}(w_{y_i}; D_i^{\text{train}}) + \rho \cdot \Omega(W_x, w_{y_i})$

Cai, H., et al. *Once for All: Train One Network and Specialize it for Efficient Deployment*. ICLR 2020.

Learn a shared model of the environment allowing robust performance across all agents.

	Problem mapping
$i \in [n]$	Different agents in the environment
UL var x	Environment model params E_x
LL var y_i	Per-agent params A_{y_i}
UL obj $\mathbb{E}_{\xi_i} f_i(x, y_i; \xi_i)$	$-\mathcal{R}(A_{y_i}, E_x; R^{\text{val}})$
LL obj $\mathbb{E}_{\zeta_i} g_i(x, y_i; \zeta_i)$	$-\mathcal{R}(A_{y_i}, E_x; R^{\text{train}})$

Flexibility: Diff agents can have diff action spaces (land, air, water).

$$\min_{x \in \mathcal{X} \subseteq \mathbb{R}^D} \max_{i \in [n]} F_i(x) \quad \text{subject to} \quad y_i^*(x) \in \arg \min_{y_i \in \mathcal{Y}_i = \mathbb{R}^{d_i}} G_i(x, y_i) \quad \forall i \in [n]$$

Reformulation

$$\min_{x \in \mathcal{X} \subseteq \mathbb{R}^D} \max_{\lambda \in \Delta_n} \sum_{i=1}^n \lambda_i F_i(x) \quad \text{s.t.} \quad y_i^*(x) \in \arg \min_{y_i \in \mathcal{Y}_i = \mathbb{R}^{d_i}} G_i(x, y_i) \quad \forall i \in [n]$$

$$\Delta_n = \left\{ \lambda \in \mathbb{R}^n : \lambda_i \geq 0 \quad \forall i \in [n], \sum_{i=1}^n \lambda_i = 1 \right\} \quad (\text{n-simplex})$$

MORBiT: Multi-Objective Robust Bilevel Two-timescale alg

Algorithm 1: Learning rates α, β, γ for x, y, λ resp

for $k = 1, 2, \dots, K$ **do**

$$\left[\begin{array}{ll} \mathbf{y}^{(k+1)} \leftarrow \mathbf{y}^{(k)} - \beta \mathbf{h}^{(k)} & \text{(SGD)} \\ x^{(k+1)} \leftarrow \text{proj}_{\mathcal{X}}(x^{(k)} - \alpha h_x^{(k)}) & \text{(Proj SGD)} \\ \lambda^{(k+1)} \leftarrow \text{proj}_{\Delta_n}(\lambda^{(k)} + \gamma h_\lambda^{(k)}) & \text{(Proj SGA)} \end{array} \right.$$

 Sample $\tau \sim \mathcal{U}(\{1, \dots, K\})$
return $\bar{x} \leftarrow x^{(\tau)}, \bar{y}_i \leftarrow y_i^{(\tau-1)}, \bar{\lambda} \leftarrow \lambda^{(\tau)}$

For $n = 1$, it reduces to the TTSA algorithm

Hong, M., et al. *A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic*. arXiv 2020.

```

for  $k = 1, 2, \dots, K$  do
     $\mathbf{y}^{(k+1)} \leftarrow \mathbf{y}^{(k)} - \beta \mathbf{h}^{(k)}$       (SGD)
     $\mathbf{x}^{(k+1)} \leftarrow \text{proj}_{\mathcal{X}}(\mathbf{x}^{(k)} - \alpha \mathbf{h}_x^{(k)})$   (P-SGD)
     $\lambda^{(k+1)} \leftarrow \text{proj}_{\Delta_n}(\lambda^{(k)} + \gamma \mathbf{h}_\lambda^{(k)})$  (P-SGA)
Sample  $\tau \sim \mathcal{U}(\{1, \dots, K\})$ 
return  $\bar{\mathbf{x}} \leftarrow \mathbf{x}^{(\tau)}, \bar{\mathbf{y}}_i \leftarrow \mathbf{y}_i^{(\tau-1)}, \bar{\lambda} \leftarrow \lambda^{(\tau)}$ 
    
```

- ▶ $\mathbf{y}^{(k)} \triangleq \{\mathbf{y}_1^{(k)}, \dots, \mathbf{y}_n^{(k)}\}$
- ▶ $\mathbf{h}^{(k)} \triangleq \{\nabla_{\mathbf{y}_1} g_1(\mathbf{x}^{(k)}, \mathbf{y}_1^{(k)}; \zeta_1), \dots, \nabla_{\mathbf{y}_n} g_n(\mathbf{x}^{(k)}, \mathbf{y}_n^{(k)}; \zeta_n)\}$
- ▶ $\mathbf{h}_x^{(k)} \approx \sum_{i \in [n]} \lambda_i^{(k)} \underbrace{\bar{\nabla} f_i(\mathbf{x}^{(k)}, \mathbf{y}_i^{(k+1)}; \xi_i)}_{**}$

$$** \triangleq \nabla_{\mathbf{x}} f_i(\mathbf{x}, \mathbf{y}_i; \xi_i) - \nabla_{\mathbf{x} \mathbf{y}_i}^2 g_i(\mathbf{x}, \mathbf{y}_i; \zeta_i) [\nabla_{\mathbf{y}_i}^2 g_i(\mathbf{x}, \mathbf{y}_i; \zeta_i)]^{-1} \nabla_{\mathbf{y}_i} f_i(\mathbf{x}, \mathbf{y}_i; \xi_i)$$

- ▶ $\mathbf{h}_\lambda^{(k)} \triangleq [f_1(\mathbf{x}^{(k)}, \mathbf{y}_1^{(k+1)}; \xi_1), \dots, f_n(\mathbf{x}^{(k)}, \mathbf{y}_n^{(k+1)}; \xi_n)]^\top$

Learning rates

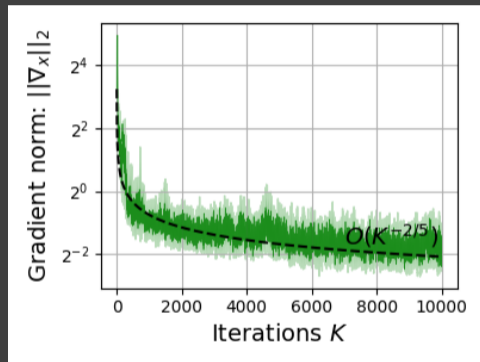
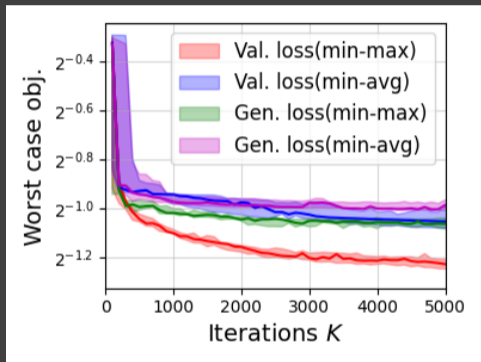
$$\alpha \sim \mathcal{O}\left(K^{-\frac{3}{5}}\right), \beta \sim \mathcal{O}\left(K^{-\frac{2}{5}}\right), \gamma \sim \mathcal{O}\left(\sqrt{n}K^{-\frac{3}{5}}\right)$$

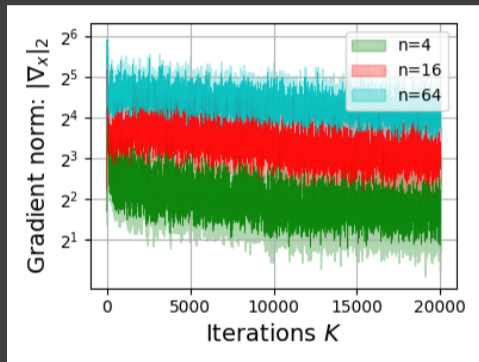
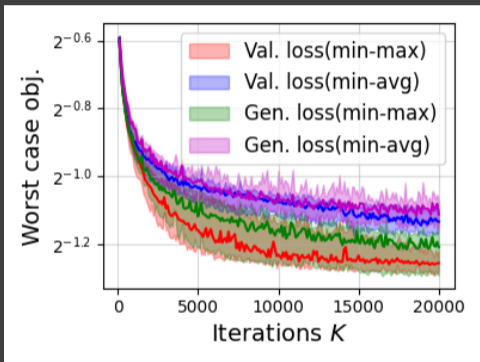
Convergence rate of MORBiT

$$\begin{aligned} \mathbb{E}[\|\hat{\mathbf{x}}(\bar{\mathbf{x}}) - \bar{\mathbf{x}}\|^2] &\leq \tilde{\mathcal{O}}(\sqrt{n}K^{-2/5})^\dagger, \\ \mathbb{E}\left[\max_{i \in [n]} \|\bar{\mathbf{y}}_i - \mathbf{y}_i^*(\bar{\mathbf{x}})\|^2\right] &\leq \tilde{\mathcal{O}}(\sqrt{n}K^{-2/5}), \\ \max_{\lambda} \mathbb{E}[F(\bar{\mathbf{x}}, \lambda)] - \mathbb{E}[F(\bar{\mathbf{x}}, \bar{\lambda})] &\leq \tilde{\mathcal{O}}(\sqrt{n}K^{-2/5}). \end{aligned}$$

$^\dagger \hat{\mathbf{x}}(\bar{\mathbf{x}})$ is the proximal map

- ▶ Handling the non-smooth $\max_{i \in [n]}$
- ▶ Establishing convergence of each \bar{y}_i *simultaneously* $\forall i \in [n]$
- ▶ Descent equation involving $(n + 1)$ sequences
- ▶ Establishing convergence of λ





Caveats:

- ▶ Useful when gap between max and mean is large
- ▶ From a learning perspective, \mathcal{X} should have enough capacity to simultaneously optimize all UL objectives
- ▶ Room for improvement for large n

Arxiv Preprint: <https://arxiv.org/pdf/2203.01924.pdf>

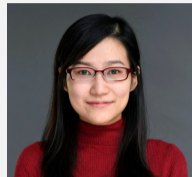
Gu, Alex, et al. *Min-Max Bilevel Multi-objective Optimization with Applications in Machine Learning*. arXiv 2022.



Alex Gu



Songtao Lu



Tsui-Wei (Lily)
Weng

Thank You!

Problem/Method	Bilevel	Multi-obj	Min-max	Single-loop	$\mathcal{X} \subset \mathbb{R}^{d_x}$	$y_i \subset \mathbb{R}^{d_{y_i}}$
Distributionally Robust Learning	†	✗	✓	-	✓	✗
Adversarially Robust Learning	†	✗	✓	-	✗	✓
Multi-task Learning (MTL)	†	□	✗	-	✗	✗
Robust MTL [Mehta et al., 2012]	†	✓	✓	-	✗	✗
Meta-learning	†	□	✗	-	✗	✗
HiBSA [Lu et al., 2020]	✗	✗	✓	✓	✓	-
GDA [Lin et al., 2020]	✗	✗	✓	✓	✗	-
TR-MAML [Collins et al., 2020]	✗	✓	✓	✓	✓	-
BSA [Ghadimi and Wang, 2018]	✓	✗	✗	✗	✓	✗
TTSA [Hong et al., 2020]	✓	✗	✗	✓	✓	✗
StocBio [Ji et al., 2021]	✓	✗	✗	✗	✗	✗
MRBO [Yang et al., 2021]	✓	✗	✗	✓	✗	✗
VRBO [Yang et al., 2021]	✓	✗	✗	✗	✗	✗
ALSET [Chen et al., 2021]	✓	✗	✗	✓	✗	✗
STABLE [Chen et al., 2022]	✓	✗	✗	✓	✓	✗
MMB [Hu et al., 2022]	✓	✗	✓	✓	✗	✗
MORBiT (Ours)	✓	✓	✓	✓	✓	✗

- Nishant A. Mehta, Dongryeol Lee, and Alexander G. Gray. Minimax multi-task learning and a generalized loss-compositional paradigm for mtl. In *Proceedings of the 25th International Conference on Neural Information Processing Systems*, page 2150–2158, 2012.
- Songtao Lu, Ioannis Tsaknakis, Mingyi Hong, and Yongxin Chen. Hybrid block successive approximation for one-sided non-convex min-max problems: algorithms and applications. *IEEE Transactions on Signal Processing*, 68:3676–3691, 2020.
- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *Proceedings of International Conference on Machine Learning (ICML)*, pages 6083–6093. PMLR, 2020.
- Liam Collins, Aryan Mokhtari, and Sanjay Shakkottai. Task-robust model-agnostic meta-learning. *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020.
- Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.
- Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale framework for bilevel optimization: Complexity analysis and application to actor-critic. *arXiv preprint arXiv:2007.05170*, 2020. URL <https://arxiv.org/pdf/2007.05170.pdf>.
- Kaiyi Ji, Junjie Yang, and Yingbin Liang. Bilevel optimization: Convergence analysis and enhanced design. In *Proceedings of International Conference on Machine Learning (ICML)*, pages 4882–4892. PMLR, 2021.
- Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker participation in non-iid federated learning. *arXiv preprint arXiv:2101.11203*, 2021.
- Tianyi Chen, Yuejiao Sun, and Wotao Yin. Closing the gap: Tighter analysis of alternating stochastic gradient methods for bilevel problems. *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 34:25294–25307, 2021.
- Tianyi Chen, Yuejiao Sun, Quan Xiao, and Wotao Yin. A single-timescale method for stochastic bilevel optimization. In *Proceedings of International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 2466–2488. PMLR, 2022.
- Quanqi Hu, Yongjian Zhong, and Tianbao Yang. Multi-block min-max bilevel optimization with applications in multi-task deep auc maximization. *arXiv preprint arXiv:2206.00260*, 2022.