

Overcoming Lower-Level Constraints in Bilevel Optimization: A Novel Approach with Regularized Gap Functions

Wei Yao, Haian Yin, Shangzhi Zeng, Jin Zhang

Problem Formulation

The **constrained** bilevel optimization (BiO) is formulated as:

$$\min_{x \in X, y \in Y} F(x, y) \quad \text{s.t.} \quad y \in S(x) := \operatorname{argmin}_{y \in Y} \left\{ f(x, y) \text{ s.t. } g(x, y) \leq 0 \right\},$$

which tackles nested structures present in constrained learning tasks like constrained meta-learning, adversarial learning, and distributed bilevel optimization.

Assumptions:

- F is **Smooth** and **bounded below** over the feasible set.
- f is **Convex** in y , and **smooth** in both variables.
- g is **Convex** in y , **smooth**, and **Lipschitz continuous** (including its gradients).

Goal: develop a single-loop, first-order algorithm without projection onto the coupled lower-level constraint set.

Gap Function and Its Property

Here we introduce the doubly regularized gap function for the lower-level problem:

$$\mathcal{G}_\gamma(x, y, z) := \max_{\theta \in Y, \lambda \in \mathbb{R}_+^p} \left\{ \mathcal{L}(x, y, \lambda) - \frac{1}{2\gamma_2} \|\lambda - z\|^2 - \mathcal{L}(x, \theta, z) - \frac{1}{2\gamma_1} \|\theta - y\|^2 \right\},$$

where the Lagrangian function $\mathcal{L}(x, y, z) := f(x, y) + z^T g(x, y)$.

Property I (Lemma 2.1.)

$$\mathcal{G}_\gamma(x, y, z) \leq 0 \Leftrightarrow y \in S(x) \text{ and } z \in \mathcal{M}(x, y)$$

Property II (Lemma 2.2.)

$$\nabla \mathcal{G}_\gamma(x, y, z) = \begin{bmatrix} \nabla_x f(x, y) + (\lambda^*)^T \nabla_x g(x, y) \\ \nabla_y f(x, y) + (\lambda^*)^T \nabla_y g(x, y) \\ -(z - \lambda^*) / \gamma_2 \end{bmatrix} - \begin{bmatrix} \nabla_x f(x, \theta^*) + z^T \nabla_x g(x, \theta^*) \\ (y - \theta^*) / \gamma_1 \\ g(x, \theta^*) \end{bmatrix}$$

$$\theta^* := \theta^*(x, y, z) := \operatorname{argmin}_{\theta \in Y} \left\{ \mathcal{L}(x, \theta, z) + \frac{1}{2\gamma_1} \|\theta - y\|^2 \right\},$$

$$\lambda^* := \lambda^*(x, y, z) := \operatorname{argmax}_{\lambda \in \mathbb{R}_+^p} \left\{ \mathcal{L}(x, y, \lambda) - \frac{1}{2\gamma_2} \|\lambda - z\|^2 \right\} = \operatorname{Proj}_{\mathbb{R}_+^p} (z + \gamma_2 g(x, y)).$$

Reformulation

Based on the Property I of gap function, we can reformulate the BiO **equivalently** as

$$\min_{(x,y,z) \in X \times Y \times \mathbb{R}_+^p} F(x, y) \quad \text{s.t.} \quad \mathcal{G}_\gamma(x, y, z) \leq 0.$$

To develop a gradient-based algorithm, we explore its **penalty** formulation:

$$\min_{(x,y,z) \in X \times Y \times Z} F(x, y) + c \mathcal{G}_\gamma(x, y, z),$$

Main Algorithm

➤ Bilevel Constrained GAp Function-based First-order Algorithm (BiC-GAFFA)

1. Update Auxiliary Variables:

$$\theta^{k+1} = \text{Proj}_Y \left(\theta^k - \eta_k \nabla_\theta (\mathcal{L}(x^k, \theta^k, z^k) + \frac{1}{2\gamma_1} \|\theta^k - y^k\|^2) \right)$$

$$\lambda^{k+1} = \text{Proj}_{\mathbb{R}_+^p} (z^k + \gamma_2 g(x^k, y^k))$$

2. Update Main Variables:

$$\begin{aligned} & (x^{k+1}, y^{k+1}, z^{k+1}) \\ &= \text{Proj}_{X \times Y \times Z} \left((x^k, y^k, z^k) - \alpha_k \nabla \left(\frac{1}{c_k} F(x^k, y^k) + \mathcal{G}_\gamma(x^k, y^k, z^k) \right) \right) \end{aligned}$$

Extension to Bilevel Optimization with Minimax Lower-level Problem

For bilevel optimization with minimax lower-level problem:

$$\min_{x \in X, y \in Y, z \in Z} F(x, y, z) \quad \text{s.t.} \quad (y, z) \in \mathcal{SP}(x) := \text{Sol} \left[\min_{y \in Y} \max_{z \in Z} f(x, y, z) \right].$$

We use following gap function:

$$\mathcal{G}_\gamma^{\text{saddle}}(x, y, z) := \max_{\theta \in Y, \lambda \in \mathbb{R}_+^p} \left\{ f(x, y, z) - \frac{1}{2\gamma_2} \|\lambda - z\|^2 - f(x, \theta, z) - \frac{1}{2\gamma_1} \|\theta - y\|^2 \right\}.$$

And reformulate the BiO as:

$$\min_{(x, y, z) \in X \times Y \times \mathbb{R}_+^p} F(x, y, z) \quad \text{s.t.} \quad \mathcal{G}_\gamma^{\text{saddle}}(x, y, z) \leq 0.$$

➤ Single-loop Hessian-free algorithm for BiO with minimax lower-level problem

1. Update Auxiliary Variables:

$$\theta^{k+1} = \text{Proj}_Y \left(\theta^k - \eta_k \nabla_\theta (\mathcal{L}(x^k, \theta^k, z^k) + \frac{1}{2\gamma_1} \|\theta^k - y^k\|^2) \right)$$

$$\lambda^{k+1} = \text{Proj}_{\mathbb{R}_+^p} (z^k + \gamma_2 g(x^k, y^k))$$

2. Update Main Variables:

$$\begin{aligned} & (x^{k+1}, y^{k+1}, z^{k+1}) \\ &= \text{Proj}_{X \times Y \times Z} \left((x^k, y^k, z^k) - \alpha_k \nabla \left(\frac{1}{c_k} F(x^k, y^k) + \tilde{\mathcal{G}}_\gamma^{\text{saddle}}(x^k, y^k, z^k) \right) \right) \end{aligned}$$

Synthetic Experiments

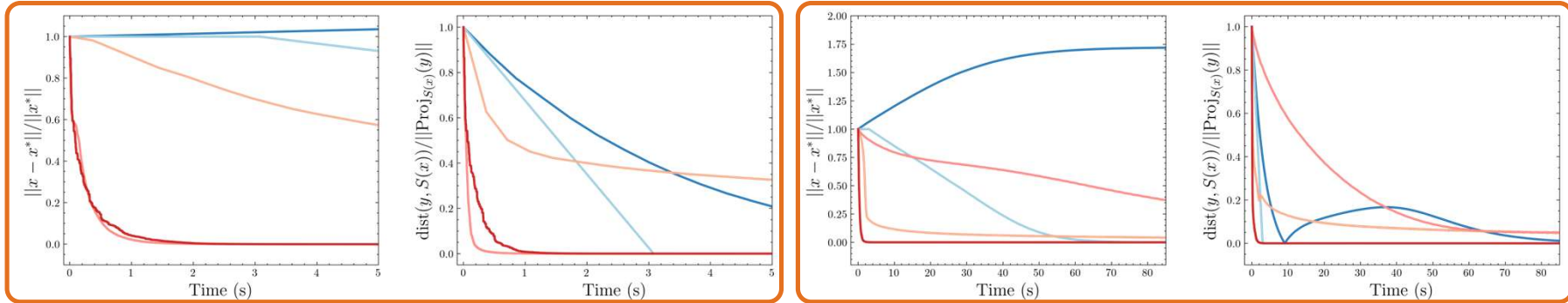
$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{X}, (\mathbf{y}_1, \mathbf{y}_2) \in \mathcal{Y}} \quad & (\mathbf{y}_1 - 2 \cdot \mathbf{1}_n)^T (\mathbf{x} - \mathbf{1}_n) + \|\mathbf{y}_2 + 3 \cdot \mathbf{1}_n\|^2 \\ \text{s.t.} \quad & (\mathbf{y}_1, \mathbf{y}_2) \in \underset{(\mathbf{y}_1, \mathbf{y}_2) \in \mathcal{Y}}{\operatorname{argmin}} \left\{ \frac{1}{2} \|\mathbf{y}_1\|^2 - \mathbf{x}^T \mathbf{y}_1 + \mathbf{1}_n^T \mathbf{y}_2 \mid \sum_{i=1}^n h(\mathbf{x}_i) + \mathbf{1}_n^T \mathbf{y}_1 + \mathbf{1}_n^T \mathbf{y}_2 = 0 \right\}. \end{aligned}$$

$$\textcircled{1} \ h(x) = x$$

$$\textcircled{2} \ h(x) = x^3$$

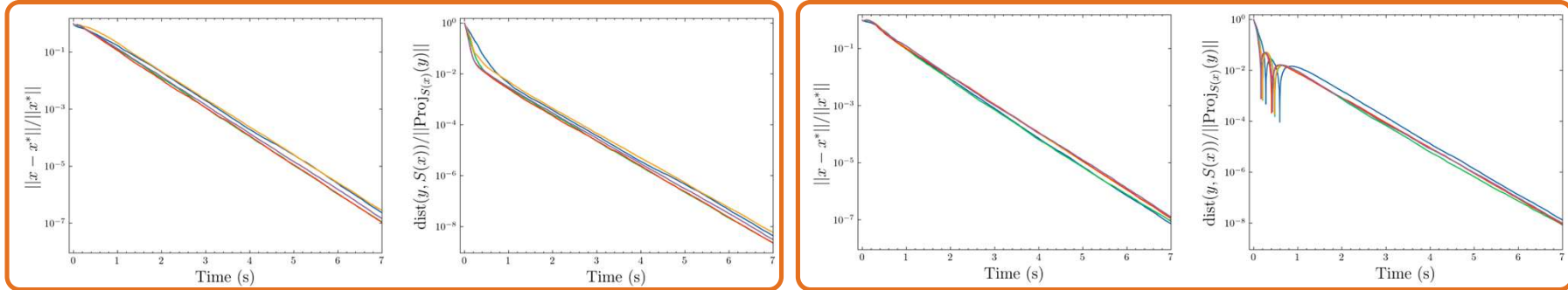
Performance (Compared with the SOTA Algorithms)

— AIPOD — GAM — BVFSM — LV-HBA — BiC-GAFFA



Sensitivity Analysis on Parameters

— $\gamma_1=1.0$ — $\gamma_1=3.0$ — $\gamma_1=5.0$ — $\gamma_1=7.0$ — $\gamma_1=10.0$



Hyperparameter Optimization on Sparse Group Lasso Problem

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p, \boldsymbol{\lambda} \in \mathbb{R}_+^{M+1}} \frac{1}{2} \sum_{i \in I_{\text{val}}} |b_i - \boldsymbol{\beta}^T \mathbf{a}_i|^2$$

$$\text{s.t.} \quad \boldsymbol{\beta} \in \arg \min_{\hat{\boldsymbol{\beta}} \in \mathbb{R}^p} \left\{ \frac{1}{2} \sum_{i \in I_{\text{tr}}} |b_i - \hat{\boldsymbol{\beta}}^T \mathbf{a}_i|^2 + \sum_{m=1}^M \lambda_m \|\hat{\boldsymbol{\beta}}^{(m)}\|_2 + \lambda_{M+1} \|\hat{\boldsymbol{\beta}}\|_1 \right\}.$$

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^p, \mathbf{u} \in \mathbb{R}_+^{M+1}} \frac{1}{2} \sum_{i \in I_{\text{val}}} |b_i - \boldsymbol{\beta}^T \mathbf{a}_i|^2$$

$$\text{s.t.} \quad \boldsymbol{\beta} \in \arg \min_{\hat{\boldsymbol{\beta}} \in \mathbb{R}^p} \left\{ \frac{1}{2} \sum_{i \in I_{\text{tr}}} |b_i - \hat{\boldsymbol{\beta}}^T \mathbf{a}_i|^2 \text{ s.t. } \|\hat{\boldsymbol{\beta}}^{(m)}\|_2^2 \leq u_m, m \in [M], \|\hat{\boldsymbol{\beta}}\|_1 \leq u_{M+1} \right\}.$$

Hyperparameter

Decoupling

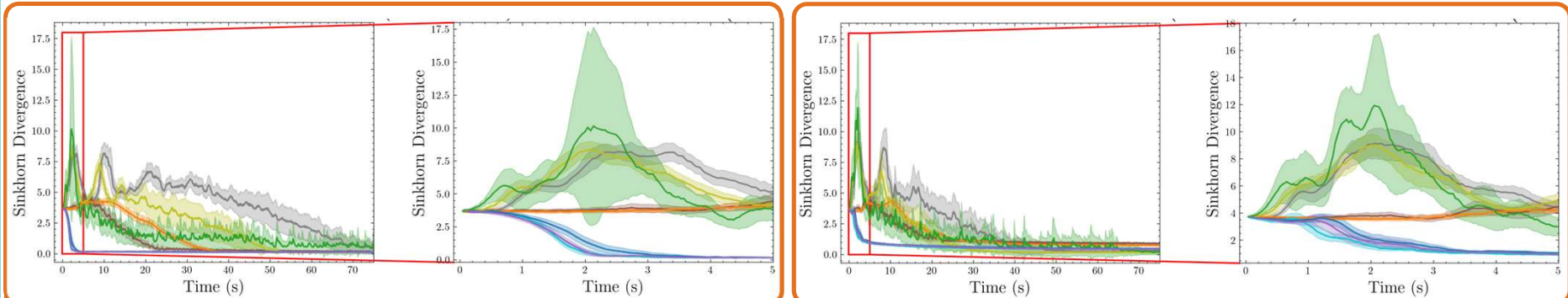
Method	nTr = 100, nVal = 100, nTest = 300			nTr = 300, nVal = 300, nTest = 300		
	Time (s)	Val Err	Test Err	Time (s)	Val Err	Test Err
Grid	17.3 ± 0.9	35.9 ± 7.2	37.7 ± 6.7	78.7 ± 1.9	18.9 ± 2.3	19.8 ± 1.8
Random	17.4 ± 0.7	33.6 ± 6.7	35.7 ± 6.2	78.6 ± 2.5	18.7 ± 2.4	19.5 ± 1.9
TPE	16.9 ± 0.7	33.9 ± 7.0	36.0 ± 5.6	74.7 ± 2.2	18.9 ± 2.3	19.8 ± 1.9
IGJO	21.2 ± 2.2	19.7 ± 2.8	25.6 ± 4.4	49.9 ± 2.6	16.5 ± 2.5	18.1 ± 1.4
VF-iDCA	12.4 ± 0.5	14.6 ± 2.6	25.4 ± 3.9	40.7 ± 1.7	14.9 ± 2.1	17.2 ± 1.3
BiC-GAFFA	21.4 ± 0.7	7.3 ± 1.3	22.3 ± 3.0	22.0 ± 1.0	12.8 ± 1.4	17.1 ± 1.3

GAN with or with out Constraint

Model: $\min_{G,D} \mathcal{L}_{\text{gen}}(G, D) \quad \text{s.t.} \quad D \in \operatorname{argmin}_{D \in \mathcal{D}} \mathcal{L}_{\text{det}}(G, D).$

WGAN: $\mathcal{L}_{\text{gen}}(G, D) = -\mathbb{E}_{\mathbf{z} \sim \mathbb{N}}[D(G(\mathbf{z}))],$
 $\mathcal{L}_{\text{det}}(G, D) = -\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r}[D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim \mathbb{N}}[D(G(\mathbf{z}))],$
 $\mathcal{D} := \left\{ D(\mathbf{x}) \quad \text{s.t.} \quad \max_{\hat{\mathbf{x}} \in \mathbb{P}_{\hat{\mathbf{x}}}} \|\nabla D(\hat{\mathbf{x}})\|_2 \leq 1 \right\}.$

— GAN — WGAN-GP — UGAN — Bi-WGAN (BiC-GAFFA)
— WGAN — Con-GAN — Bi-GAN (BiC-GAFFA) — Bi-ConGAN (BiC-GAFFA)



Thanks

Overcoming Lower-Level Constraints in Bilevel Optimization: A Novel Approach with Regularized Gap Functions

Wei Yao^{1,2}, Haian Yin², Shangzhi Zeng^{1,2}, Jin Zhang^{2,1,3} Emails: {yaow, yinha, zengsz, zhangj9}@sustech.edu.cn
¹National Center for Applied Mathematics Shenzhen, ²Southern University of Science and Technology, ³Smart City Research Institute of CETC

Thirteenth International Conference
on Learning Representations

Problem Formulation

The **constrained bilevel optimization (BiO)** is formulated as:

$$\min_{x \in X, y \in Y} F(x, y) \quad \text{s.t.} \quad y \in S(x) := \arg\min_{y \in Y} \{f(x, y) \mid g(x, y) \leq 0\},$$

which tackles nested structures present in constrained learning tasks like constrained meta-learning, adversarial learning, and distributed bilevel optimization.

Assumptions:

- F is **Smooth** and **bounded below** over the feasible set.
- f is **Convex** in y , and **smooth** in both variables.
- g is **Convex** in y , **smooth**, and **Lipschitz continuous** (including its gradients).

Goal: develop a single-loop, first-order algorithm without projection onto the coupled lower-level constraint set.

Main Algorithm

➤ Bilevel Constrained Gap Function-based First-order Algorithm (BiC-GAFFA)

- Update Auxiliary Variables:
 $\theta^{k+1} = \text{Proj}_Y \left(\theta^k - \eta_k \nabla_\theta (\mathcal{L}(x^k, \theta^k, z^k) + \frac{1}{2\gamma_1} \|\theta^k - y^k\|^2) \right)$
 $\lambda^{k+1} = \text{Proj}_{\mathbb{R}_+^n} (z^k + \gamma_2 g(x^k, y^k))$
- Update Main Variables:
 $(x^{k+1}, y^{k+1}, z^{k+1}) = \text{Proj}_{X \times Y \times Z} \left((x^k, y^k, z^k) - \alpha_k \nabla \left(\frac{1}{c_k} F(x^k, y^k) + \mathcal{G}_\gamma(x^k, y^k, z^k) \right) \right)$

Hyperparameter Optimization on Sparse Group Lasso Problem

$$\min_{\beta \in \mathbb{R}^p, \lambda \in \mathbb{R}_+^{M+1}} \frac{1}{2} \sum_{i \in I_{\text{out}}} |b_i - \beta^T \mathbf{a}_i|^2$$

$$\text{s.t.} \quad \beta \in \arg\min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \sum_{i \in I_{\text{in}}} |b_i - \beta^T \mathbf{a}_i|^2 + \sum_{m=1}^M \lambda_m \|\beta^{(m)}\|_2 + \lambda_{M+1} \|\beta\|_1 \right\}.$$

Hyperparameter Decoupling

$$\text{s.t.} \quad \beta \in \arg\min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \sum_{i \in I_{\text{in}}} |b_i - \beta^T \mathbf{a}_i|^2 \mid \|\beta^{(m)}\|_2^2 \leq u_m, m \in [M], \|\beta\|_1 \leq u_{M+1} \right\}.$$

Method	nTr = 100, nVal = 100, nTest = 300			nTr = 300, nVal = 300, nTest = 300		
	Time (s)	Val Err	Test Err	Time (s)	Val Err	Test Err
Grid	17.3 ± 0.9	35.9 ± 7.2	37.7 ± 6.7	78.7 ± 1.9	18.9 ± 2.3	19.8 ± 1.8
Random	17.4 ± 0.7	33.6 ± 6.7	35.7 ± 6.2	78.6 ± 2.5	18.7 ± 2.4	19.5 ± 1.9
TPE	16.9 ± 0.7	33.9 ± 7.0	36.0 ± 5.6	74.7 ± 2.2	18.9 ± 2.3	19.8 ± 1.9
IGJO	21.2 ± 2.2	19.7 ± 2.8	25.6 ± 4.4	49.9 ± 2.6	16.5 ± 2.5	18.1 ± 1.4
VF-IDCA	12.4 ± 0.5	14.6 ± 2.6	25.4 ± 3.9	40.7 ± 1.7	14.9 ± 2.1	17.2 ± 1.3
BiC-GAFFA	21.4 ± 0.7	7.3 ± 1.3	22.3 ± 3.0	22.0 ± 1.0	12.8 ± 1.4	17.1 ± 1.3

Gap Function and Its Property

Here we introduce the doubly regularized gap function for the lower-level problem:

$$\mathcal{G}_\gamma(x, y, z) := \max_{\theta \in Y, \lambda \in \mathbb{R}_+^n} \left\{ \mathcal{L}(x, y, \lambda) - \frac{1}{2\gamma_2} \|\lambda - z\|^2 - \mathcal{L}(x, \theta, z) - \frac{1}{2\gamma_1} \|\theta - y\|^2 \right\},$$

where the Lagrangian function $\mathcal{L}(x, y, z) := f(x, y) + z^T g(x, y)$.

Property I (Lemma 2.1.)
 $\mathcal{G}_\gamma(x, y, z) \leq 0 \Leftrightarrow y \in S(x)$ and $z \in \mathcal{M}(x, y)$

Property II (Lemma 2.2.)

$$\nabla \mathcal{G}_\gamma(x, y, z) = \begin{bmatrix} \nabla_x f(x, y) + (\lambda^*)^T \nabla_x g(x, y) \\ \nabla_y f(x, y) + (\lambda^*)^T \nabla_y g(x, y) \\ -(z - \lambda^*) / \gamma_2 \end{bmatrix} - \begin{bmatrix} \nabla_x f(x, \theta^*) + z^T \nabla_x g(x, \theta^*) \\ (y - \theta^*) / \gamma_1 \\ g(x, \theta^*) \end{bmatrix}$$

$$\theta^* := \theta^*(x, y, z) := \arg\min_{\theta \in Y} \left\{ \mathcal{L}(x, \theta, z) + \frac{1}{2\gamma_1} \|\theta - y\|^2 \right\},$$

$$\lambda^* := \lambda^*(x, y, z) := \arg\max_{\lambda \in \mathbb{R}_+^n} \left\{ \mathcal{L}(x, y, z) - \frac{1}{2\gamma_2} \|\lambda - z\|^2 \right\} = \text{Proj}_{\mathbb{R}_+^n} (z + \gamma_2 g(x, y)).$$

Extension to Bilevel Optimization with Minimax Lower-level Problem

For bilevel optimization with minimax lower-level problem:

$$\min_{x \in X, y \in Y, z \in Z} F(x, y, z) \quad \text{s.t.} \quad (y, z) \in \mathcal{SP}(x) := \text{Sol} \left[\min_{y \in Y} \max_{z \in Z} f(x, y, z) \right].$$

We use following gap function:

$$\mathcal{G}_\gamma^{\text{middle}}(x, y, z) := \max_{\theta \in Y, \lambda \in \mathbb{R}_+^n} \left\{ f(x, y, z) - \frac{1}{2\gamma_2} \|\lambda - z\|^2 - f(x, \theta, z) - \frac{1}{2\gamma_1} \|\theta - y\|^2 \right\}.$$

And reformulate the BiO as:

$$\min_{(x, y, z) \in X \times Y \times \mathbb{R}_+^n} F(x, y, z) \quad \text{s.t.} \quad \mathcal{G}_\gamma^{\text{middle}}(x, y, z) \leq 0.$$

➤ Single-loop Hessian-free algorithm for BiO with minimax lower-level problem

- Update Auxiliary Variables:
 $\theta^{k+1} = \text{Proj}_Y \left(\theta^k - \eta_k \nabla_\theta (\mathcal{L}(x^k, \theta^k, z^k) + \frac{1}{2\gamma_1} \|\theta^k - y^k\|^2) \right)$
 $\lambda^{k+1} = \text{Proj}_{\mathbb{R}_+^n} (z^k + \gamma_2 g(x^k, y^k))$
- Update Main Variables:
 $(x^{k+1}, y^{k+1}, z^{k+1}) = \text{Proj}_{X \times Y \times Z} \left((x^k, y^k, z^k) - \alpha_k \nabla \left(\frac{1}{c_k} F(x^k, y^k) + \tilde{\mathcal{G}}_\gamma^{\text{middle}}(x^k, y^k, z^k) \right) \right)$

Hyper-Cleaning on SVM

$$\min_{\mathbf{c}, \mathbf{w}, b, \xi} \text{Logistic}_{D_{\text{out}}}(\mathbf{w}, b) \quad \text{s.t.} \quad (\mathbf{w}, b, \xi) \in \text{SVM}_{D_{\text{in}}}(\mathbf{w}, b, \xi; \mathbf{c})$$

Reformulation

Based on the Property I of gap function, we can reformulate the BiO **equivalently** as

$$\min_{(x, y, z) \in X \times Y \times \mathbb{R}_+^n} F(x, y) \quad \text{s.t.} \quad \mathcal{G}_\gamma(x, y, z) \leq 0.$$

To develop a gradient-based algorithm, we explore its **penalty** formulation:

$$\min_{(x, y, z) \in X \times Y \times \mathbb{R}_+^n} F(x, y) + c \mathcal{G}_\gamma(x, y, z),$$

Synthetic Experiments

GAN with or without Constraint

$$\min_{G, D} \mathcal{L}_{\text{gen}}(G, D) \quad \text{s.t.} \quad D \in \arg\min_{D \in \mathcal{D}} \mathcal{L}_{\text{det}}(G, D).$$