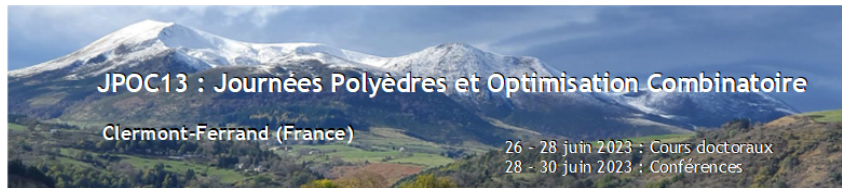


Stochastic Bi-level Problems: Variants, Algorithms, and Implementation

Combinatorial Optimization and Machine Learning | Lecture 6

Sanjeeb Dash / Parikshit Ram

June 27, 2023



1 Bi-level Problem Variants & Algorithms

- Handling Constraints
- Non-Singleton LL
- Specialized solvers
- Multi-objective Bi-level

2 Implementation Details

3 Further Reading

1 Bi-level Problem Variants & Algorithms

- Handling Constraints
- Non-Singleton LL
- Specialized solvers
- Multi-objective Bi-level

2 Implementation Details

3 Further Reading

UL constrained $\theta \in \Theta \subset \mathbb{R}^{d_u}$ + LL unconstrained, unique solution

$$\min_{\theta \in \Theta \subset \mathbb{R}^{d_u}} f_u(\theta, \phi^*(\theta)) \quad \text{subject to} \quad \phi^*(\theta) = \arg \min_{\phi \in \mathbb{R}^{d_l}} f_l(\theta, \phi) \quad (1)$$

Handle UL constrained via **projected gradient descent** for UL update (TTSA [Hong et al., 2020, 2023], STABLE [Chen et al., 2022])

UL constrained + LL constrained $\phi \in \Phi(\theta) \subset \mathbb{R}^{d_l}$, unique solution

$$\min_{\theta \in \Theta \subset \mathbb{R}^{d_u}} f_u(\theta, \phi^*(\theta)) \quad \text{subject to} \quad \phi^*(\theta) = \arg \min_{\phi \in \Phi(\theta) \subset \mathbb{R}^{d_l}} f_l(\theta, \phi) \quad (2)$$

We can obtain IG based UL gradient with additional assumptions [Giovannelli et al., 2021].

Main modification. How do we compute the **IG** $d\phi^*(\theta)^\top/d\theta$?

Since $\phi^*(\theta)$ is a LL solution, the stationarity condition in the unconstrained case gives us

$$\nabla_{\phi} f_l(\theta, \phi^*(\theta)) = 0 \quad (3)$$

Taking the derivative w.r.t. θ we have (by Implicit Function Theorem):

$$\nabla_{\theta\phi}^2 f_l(\theta, \phi^*(\theta)) + \frac{d\phi^*(\theta)}{d\theta} \underbrace{\nabla_{\phi}^2 f_l(\theta, \phi^*(\theta))}_{\text{Hessian } H} = 0. \quad (4)$$

Assuming the Hessian H is invertible,

$$\frac{d\phi^*(\theta)}{d\theta} = - \underbrace{\nabla_{\theta\phi}^2 f_l(\theta, \phi^*(\theta))}_{D_u \times D_l} \cdot \underbrace{\nabla_{\phi}^2 f_l(\theta, \phi^*(\theta))^{-1}}_{D_l \times D_l}. \quad (5)$$

LL feasible set

$$\Phi(\theta) = \left\{ \phi \in \mathbb{R}^{d_l} : P_l(\theta, \phi) \leq 0, Q_l(\theta, \phi) = 0 \right\} \quad (6)$$

At $\phi^*(\theta) \in \Phi(\theta)$, the following assumptions need to hold

- ⇒ The gradients of the active constraints are linearly independent (LICQ – linearly independent constraint qualification)
- ⇒ Strict complementarity condition holds
- ⇒ Sufficient second-order optimality conditions are satisfied

Define the Lagrangian function with Lagrange multipliers λ_P, λ_Q as

$$\mathcal{L}(\theta, \phi, \lambda_P, \lambda_Q) = f_l(\theta, \phi) + \lambda_P^\top P(\theta, \phi) + \lambda_Q^\top Q(\theta, \phi), \quad (7)$$

Assuming unique Lagrange multipliers $\lambda_P(\theta), \lambda_Q(\theta)$ at $\phi^*(\theta)$ under LICQ, this is the first-order KKT system for the LL problem:

$$\nabla_{\phi} f_l(\theta, \phi^*(\theta)) + \nabla_{\phi} P_l(\theta, \phi^*(\theta))^{\top} \cdot \lambda_P(\theta) + \nabla_{\phi} Q_l(\theta, \phi^*(\theta))^{\top} \cdot \lambda_Q(\theta) = 0 \quad (8)$$

$$\lambda_P(\theta) \odot P(\theta, \phi^*(\theta)) = 0 \quad (9)$$

$$Q(\theta, \phi^*(\theta)) = 0 \quad (10)$$

With $\varphi^*(\theta) = [\phi^*(\theta)^{\top}, \lambda_P(\theta)^{\top}, \lambda_Q(\theta)^{\top}]^{\top} \in \mathbb{R}^{(d_l+p_l+q_l)}$, we can rewrite the above as $G(\theta, \varphi^*(\theta)) = 0$.

Applying Implicit Function Theorem (under appropriate assumptions)

$$\nabla_{\theta} G^{\top} + \frac{d\varphi^*(\theta)^{\top}}{d\theta} \nabla_{\varphi} G^{\top} = 0 \quad \Rightarrow \quad \frac{d\varphi^*(\theta)^{\top}}{d\theta} = -\nabla_{\theta} G^{\top} \cdot \nabla_{\varphi} (G^{\top})^{-1} \quad (11)$$

$$\frac{d\varphi^*(\theta)^\top}{d\theta} = -\nabla_\theta G(\theta, \varphi^*(\theta))^\top \cdot \nabla_\varphi (G(\theta, \varphi^*(\theta))^\top)^{-1} \quad (12)$$

$$\nabla_\theta G(\theta, \varphi^*(\theta))^\top = \left[\underbrace{\nabla_{\theta\phi}^2 \mathcal{L}(\theta, \varphi^*(\theta))}_{\mathbb{R}^{d_u \times d_l}} \quad \underbrace{\nabla_\theta P_l(\theta, \phi^*(\theta))^\top \odot \lambda_P(\theta)^\top}_{\mathbb{R}^{d_u \times p_l}} \quad \underbrace{\nabla_\theta Q_l(\theta, \phi^*(\theta))^\top}_{\mathbb{R}^{d_u \times q_l}} \right] \quad (13)$$

$$\frac{d\varphi^*(\theta)^\top}{d\theta} = -\nabla_\theta G(\theta, \varphi^*(\theta))^\top \cdot \nabla_\varphi (G(\theta, \varphi^*(\theta))^\top)^{-1} \quad (14)$$

$$\nabla_\varphi (G(\theta, \varphi^*(\theta))^\top) = \begin{bmatrix} \underbrace{\nabla_\phi^2 \mathcal{L}(\theta, \varphi^*(\theta))}_{\mathbb{R}^{d_l \times d_l}} & \underbrace{\nabla_\phi P_l(\theta, \phi^*(\theta))^\top \odot \lambda_P(\theta)^\top}_{\mathbb{R}^{d_l \times p_l}} & \underbrace{\nabla_\phi Q_l(\theta, \phi^*(\theta))^\top}_{\mathbb{R}^{d_l \times q_l}} \\ \underbrace{\left(\nabla_\phi P_l(\theta, \phi^*(\theta))^\top \right)^\top}_{\mathbb{R}^{p_l \times d_l}} & \underbrace{\text{diag}(P_l(\theta, \phi^*(\theta)))}_{\mathbb{R}^{p_l \times p_l}} & \underbrace{0}_{\mathbb{R}^{p_l \times q_l}} \\ \underbrace{\left(\nabla_\phi Q_l(\theta, \phi^*(\theta))^\top \right)^\top}_{\mathbb{R}^{q_l \times d_l}} & \underbrace{0}_{\mathbb{R}^{q_l \times p_l}} & \underbrace{0}_{\mathbb{R}^{q_l \times q_l}} \end{bmatrix} \quad (15)$$

$$\frac{d\phi^*(\theta)^\top}{d\theta} = -\nabla_\theta G(\theta, \phi^*(\theta))^\top \cdot \nabla_\phi (G(\theta, \phi^*(\theta))^\top)^{-1} \quad (16)$$

$$\frac{d\phi^*(\theta)^\top}{d\theta} = \underbrace{\frac{d\phi^*(\theta)^\top}{d\theta}}_{\text{extract first } d_l \text{ columns of } d\phi^*(\theta)/d\theta} \cdot \begin{bmatrix} I_{d_l} \\ 0 \end{bmatrix} \quad (17)$$

LL optimal set is non-singleton

$$S(\theta) = \arg \min_{\phi \in \Phi(\theta) \subset \mathbb{R}^{d_l}} f_l(\theta, \phi) \quad (18)$$

Non-singleton LL

Implicit Gradient not available

Optimistic version.

$$\min_{\substack{\theta \in \Theta \subset \mathbb{R}^{d_u}, \\ \phi \in S(\theta)}} f_u(\theta, \phi) \quad \text{subject to} \quad S(\theta) = \arg \min_{\phi \in \Phi(\theta) \subset \mathbb{R}^{d_l}} f_l(\theta, \phi) \quad (19)$$

Best handled the surrogate optimal value function [Sow et al., 2022a]

Pessimistic version.

$$\min_{\theta \in \Theta \subset \mathbb{R}^{d_u}} \max_{\phi \in S(\theta)} f_u(\theta, \phi) \quad \text{subject to} \quad S(\theta) = \arg \min_{\phi \in \Phi(\theta) \subset \mathbb{R}^{d_l}} f_l(\theta, \phi) \quad (20)$$

Robustness via Pessimism

- ⇒ *min-min*: $\min_{\theta} \min_{\phi \in S(\theta)} f_u$ requires **just a single** $\phi \in S(\theta)$ to be “good”
- ⇒ *min-max*: $\min_{\theta} \max_{\phi \in S(\theta)} f_u$ requires **all** $\phi \in S(\theta)$ to be “good”
- ⇒ If $\min_{\phi \in S(\theta)}$ hard, inefficient *min-min* solution + approximations can be “bad”
 - ⇒ *min-max* does not require the $\min_{\phi \in S(\theta)}$ to be solved well
 - ⇒ Puts onus on \min_{θ} to find a “generally good” $S(\theta)$
 - ⇒ Makes solution “robust” to $\min_{\phi \in S(\theta)}$ quality
- ⇒ Cost of robustness: generally $\min_{\theta} \min_{\phi \in S(\theta)} f_l < \min_{\theta} \max_{\phi \in S(\theta)} f_l$ (for exact sols)

Pessimistic version.

$$\min_{\theta \in \Theta \subset \mathbb{R}^{d_u}} \max_{\phi \in S(\theta)} f_u(\theta, \phi) \quad \text{subject to} \quad S(\theta) = \arg \min_{\phi \in \Phi(\theta) \subset \mathbb{R}^{d_l}} f_l(\theta, \phi) \quad (21)$$

- ⇒ With multiple ML model types, $S(\theta)$ is the non-singleton set of “LL optimal” models for hyperparameter θ in **bi-level hyperparameter optimization**.
- ⇒ With “overparameterized” models, non-singleton $S(\theta)$ appear even in **bi-level representation learning** and other **bi-level adversarial training**
- ⇒ More situations exist ...

Robust Learning

In all such cases, we have easy access to some $\phi \in S(\theta)$, but $\min_{\phi \in S(\theta)} f_u$ is not feasible – it is more robust to find θ such that $\max_{\phi \in S(\theta)} f_u$ is optimized.

Not much attention in ML literature, so lots of opportunity!

Unconstrained singleton strongly convex LL

$$\min_{\theta \in \Theta \subset \mathbb{R}^{d_u}} f_l(\theta, \phi^*(\theta)) \quad \text{subject to} \quad \phi^*(\theta) = \arg \min_{\phi \in \mathbb{R}^{d_l}} f_l(\theta, \phi) \quad (22)$$

Main idea.

- ⇒ Solve LL with Sign-SGD [Bernstein et al., 2018]
- ⇒ Unrolled gradient simplifies UL descent step to GD

Gradient Unrolling:

For any general $t > 1$

$$\underbrace{\frac{d\varphi^{t+1}}{d\theta}}_{Z_{t+1} \in \mathbb{R}^{d_l \times d_u}} = \underbrace{\frac{\partial \varphi^{t+1}}{\partial \varphi^t}}_{A_{t+1} \in \mathbb{R}^{d_l \times d_l}} \cdot \underbrace{\frac{d\varphi^t}{d\theta}}_{Z_t} + \underbrace{\frac{\partial \varphi^{t+1}}{\partial \theta}}_{B_{t+1} \in \mathbb{R}^{d_l \times d_u}} \quad (23)$$

Recursively compute $Z_{t+1} = A_{t+1}Z_t + B_{t+1}$ by “unrolling” the gradient [Franceschi et al., 2017].

LL update with SignSGD:

$$\varphi^{t+1} \leftarrow \varphi^t - \beta \cdot \text{sign} \left(\nabla_{\phi} f_l(\theta^k, \varphi^t) \right) \quad (24)$$

Main simplification:

$$\frac{\partial}{\partial x} \text{sign}(x) = 0 \quad (\text{almost surely}) \quad (25)$$

This gives us

$$A_{t+1} = \frac{\partial \varphi^{t+1}}{\partial \varphi^t} = I_{d_l} \quad B_{t+1} = \frac{\partial \varphi^{t+1}}{\partial \theta} = 0_{d_l \times d_u} \quad (26)$$

Then the $Z_{t+1} = A_{t+1}Z_t + B_{t+1}$ recursion implies $Z_T = 0$.

Then the UL gradient reduces to:

$$\nabla_{\theta} F(\theta) = \nabla_{\theta} f_u(\theta, \phi^*(\theta)) + \frac{d\phi^*(\theta)^{\top}}{d\theta} \cdot \nabla_{\phi} f_u(\theta, \phi^*(\theta)) \quad (27)$$

$$= \nabla_{\theta} f_u(\theta, \phi^*(\theta)) \quad (28)$$

Algorithm 1 Sign-SGD based Bi-level Optimization [Fan et al., 2021]

Input: Initialization θ^0, ϕ^0 , initial learning rates α^0, β^0 for UL and LL resp.

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

 // Solve LL (approx.) with SignSGD for current θ^k

$\varphi^0 \leftarrow \phi^k$

for $t = 1, 2, \dots, T$ **do**

$\varphi^{t+1} \leftarrow \varphi^t - \beta^t \cdot \text{sign} \left(\nabla_{\phi} f_l(\theta, \phi) \Big|_{\theta=\theta^k, \phi=\varphi^t} \right)$

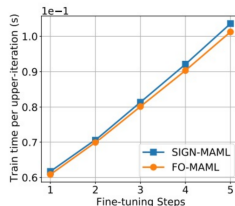
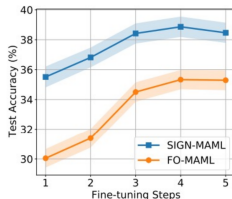
$\phi^{k+1} \leftarrow \varphi^{T+1}$

 // UL descent step with simple partial gradient

$\theta^{k+1} \leftarrow \theta^k - \alpha^k \cdot \nabla_{\theta} f_u(\theta, \phi) \Big|_{\theta=\theta^k, \phi=\phi^{k+1}}$

return θ^{K+1}, ϕ^{K+1}

- ✓ Easy to compute UL descent step (without explicitly ignoring the IG)
- ✓ Improved performance compared to methods that ignore IG, but just as fast
- ✗ SignSGD in the LL might slow convergence
- ✗ No theoretical convergence guarantees (open question!)



$$\min_{\theta \in \Theta} \{f_1(\theta), \dots, f_n(\theta)\} \quad (29)$$

- ⇒ Multiple objectives $f_i : \Theta \rightarrow \mathbb{R}$
- ⇒ Connected by single decision variable θ
- ⇒ Challenging: (usually) Trades-off between objectives
- ⇒ Dominance among solutions $\theta, \vartheta \in \Theta$:

$$\theta \succ \vartheta \Rightarrow \forall i \in [n], f_i(\theta) \leq f_i(\vartheta), \exists i \in [n], f_i(\theta) < f_i(\vartheta). \quad (30)$$

- ⇒ Pareto optimal solution θ : $\nexists \vartheta \in \Theta : \vartheta \succ \theta$
- ⇒ Pareto front: Set of (all) Pareto optimal solutions
 - ⇒ Provides a set of solutions; we might have to pick one

$$\min_{\theta \in \Theta} \{f_1(\theta), \dots, f_n(\theta)\} \quad (31)$$

Options other than Pareto optimality

⇒ Lexicographic ordering among objective – obj 1 \succ obj 2 \succ ...

⇒ A domain specific scalarization with scalars $\lambda_1, \dots, \lambda_n$:

$$\min_{\theta \in \Theta} \sum_{i \in [n]} \lambda_i f_i(\theta) \quad \text{single objective optimization} \quad (32)$$

- ⇒ We do not know the best scalarization
- ⇒ Optimize θ for the worst-case scalarization

$$\min_{\theta \in \Theta} \max_{\lambda_i, i \in [n]} \sum_{i \in [n]} \lambda_i f_i(\theta) \quad \text{minmax optimization} \quad (33)$$

$$\min_{\theta \in \Theta} \max_{\lambda_i, i \in [n]} \sum_{i \in [n]} \lambda_i f_i(\theta) \quad \text{minmax optimization} \quad (34)$$

Meaningful to put constraints on $\{\lambda_i, i \in [n]\}$

- ⇒ Non-negative $\lambda_i \geq 0 \forall i \in [n]$ – corresponds to weighted (importance) sum of objs
- ⇒ Simplex cst $\sum_{i \in [n]} \lambda_i = 1$ – weights induce a (discrete) distribution over objs
- ⇒ Gives us the n -simplex $\Delta_n = \{\lambda = [\lambda_1, \dots, \lambda_n] \in \mathbb{R}^n, \lambda_i \in [0, 1] \forall i \in [n], \mathbf{1}_n^\top \lambda = 1\}$

$$\min_{\theta \in \Theta} \max_{\substack{\lambda = [\lambda_1, \dots, \lambda_n] \\ \lambda \in \Delta_n}} \sum_{i \in [n]} \lambda_i f_i(\theta) \quad \equiv \quad \min_{\theta \in \Theta} \max_{i \in [n]} f_i(\theta) \quad (35)$$

$$\min_{\theta \in \Theta} \max_{\substack{\lambda = [\lambda_1, \dots, \lambda_n] \\ \lambda \in \Delta_n}} \sum_{i \in [n]} \lambda_i f_i(\theta) \quad \equiv \quad \min_{\theta \in \Theta} \max_{i \in [n]} f_i(\theta) \quad (36)$$

- ⇒ Optimize for the worst-case obj
- ⇒ Find θ so that all objs are “good” – robust solution
- ⇒ Worst-case scalarization – solution robust to scalarization

$$\begin{aligned} \min_{\theta \in \Theta \subset \mathbb{R}^{d_u}} \max_{i \in [n]} f_{u,i}(\theta, \phi_i^*(\theta)) &\equiv \min_{\theta \in \Theta \subset \mathbb{R}^{d_u}} \max_{\lambda \in \Delta_n} \sum_{i \in [n]} \lambda_i f_{u,i}(\theta, \phi_i^*(\theta)) \\ \text{subject to } \forall i \in [n], \phi_i^*(\theta) &= \arg \min_{\phi_i \in \Phi_i = \mathbb{R}^{d_{l,i}}} f_{l,i}(\theta, \phi_i) \end{aligned} \quad (37)$$

- ⇒ Multiple $n > 1$ UL/LL obj pairs $\{f_{u,i}, f_{l,i}\}$
- ⇒ Shared UL variable θ
- ⇒ Per-objective pair LL variable ϕ_i
- ⇒ LL variables can have different domains – that is, $\Phi_i \neq \Phi_j$
- ⇒ UL constrained
- ⇒ LL unconstrained + singleton solution

Algorithm 2 Multi-obj Robust Bi-level Two-timescale Alg [Gu et al., 2022, 2023]**Input:** Initialization $\theta^0, \lambda^0, \phi_i^0, i \in [n]$ **Input:** Initial learning rates $\alpha^0, \beta^0, \gamma^0$ for UL, LL and simplex vars resp.**for** $k = 1, 2, \dots, K$ **do**

// Single LL descent step per LL objective

$$\forall i \in [n], \phi_i^{k+1} \leftarrow \phi_i^k - \beta^k \left. \nabla_{\phi_i} f_{l,i}(\theta, \phi_i) \right|_{\theta=\theta^k, \phi_i=\phi_i^k}$$

// UL descent step with per-objective pair IG

$$\theta^{k+1} \leftarrow \mathcal{P}_{\Theta} \left(\theta^k - \alpha^k \cdot \sum_{i \in [n]} \lambda_i^k \cdot \left[\left. \nabla_{\theta} f_{u,i}(\theta, \phi_i) \right|_{\theta=\theta^k, \phi_i=\phi_i^{k+1}} - \left. \bar{\nabla} f_{u,i}(\theta, \phi_i) \right|_{\theta=\theta^k, \phi_i=\phi_i^{k+1}} \right] \right)$$

// Simplex variable ascent step

$$\lambda^{k+1} \leftarrow \mathcal{P}_{\Delta_n} \left(\lambda^k + \gamma^k \cdot [f_{u,1}(\theta^{k+1}, \phi_1^{k+1}), \dots, f_{u,n}(\theta^{k+1}, \phi_n^{k+1})]^\top \right)$$

return $\theta^{K+1}, \{\phi_i^{K+1}\}_{i \in [n]}$

Algorithm 3 Multi-obj Robust Bi-level Two-timescale Alg [Gu et al., 2022, 2023]

Input: Initialization $\theta^0, \lambda^0, \phi_i^0, i \in [n]$

Input: Initial learning rates $\alpha^0, \beta^0, \gamma^0$ for UL, LL and simplex vars resp.

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

$\forall i \in [n], \phi_i^{k+1} \leftarrow \phi_i^k - \beta^k \left[\nabla_{\phi_i} f_{l,i}(\theta, \phi_i) \right]_{\theta=\theta^k, \phi_i=\phi_i^k}$ // Single LL descent step per LL objective

$\theta^{k+1} \leftarrow \mathcal{P}_{\Theta} \left(\theta^k - \alpha^k \cdot \sum_{i \in [n]} \lambda_i^k \cdot \left[\nabla_{\theta} f_{u,i}(\theta, \phi_i) - \bar{\nabla} f_{u,i}(\theta, \phi_i) \right]_{\theta=\theta^k, \phi_i=\phi_i^{k+1}} \right)$ // UL descent step with per-objective

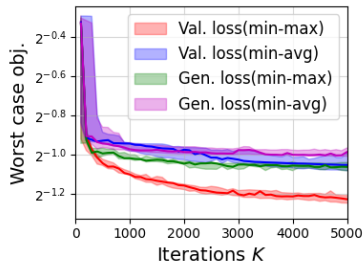
pair IG

$\lambda^{k+1} \leftarrow \mathcal{P}_{\Delta_n} \left(\lambda^k + \gamma^k \cdot [f_{u,1}(\theta^{k+1}, \phi_1^{k+1}), \dots, f_{u,n}(\theta^{k+1}, \phi_n^{k+1})]^T \right)$ // Simplex variable ascent step

return $\theta^{K+1}, \{\phi_i^{K+1}\}_{i \in [n]}$

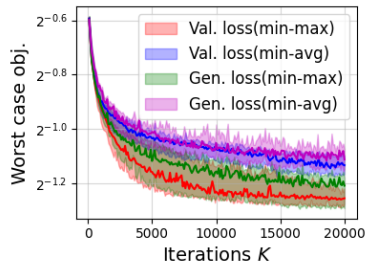
$$\bar{\nabla} f_{u,i}(\theta, \phi_i) = \nabla_{\theta, \phi_i}^2 f_{l,i}(\theta, \phi_i) \cdot \nabla_{\phi_i, \phi_i}^2 f_{l,i}(\theta, \phi_i)^{-1} \cdot \nabla_{\phi_i} f_{u,i}(\theta, \phi_i) \quad (38)$$

Representation Learning



- ⇒ UL var θ – shared representation
- ⇒ LL var ϕ_i – per-task learner

Hyperparameter Optimization



- ⇒ UL var θ – shared hyperparameter
- ⇒ LL var ϕ_i – per-task model

1 Bi-level Problem Variants & Algorithms

- Handling Constraints
- Non-Singleton LL
- Specialized solvers
- Multi-objective Bi-level

2 Implementation Details

3 Further Reading

$$\bar{\nabla} f_u(\theta, \phi) = \nabla_{\theta\phi}^2 f_l(\theta, \phi) \cdot \nabla_{\phi}^2 f_l(\theta, \phi)^{-1} \cdot \nabla_{\phi} f_u(\theta, \phi) \quad (39)$$

⇒ Inverse Hessian-vector Product (HvP) $[\nabla_{\phi}^2 f_l]^{-1} \cdot v$

⇒ Conjugate Gradient needs HvP $\nabla_{\phi}^2 f_l \cdot v$

⇒ Neumann series approx needs HvP:

$$[I - \nabla_{\phi}^2 f_l]^i \cdot v = [I - \nabla_{\phi}^2 f_l]^{i-1} \cdot [I \cdot v - \nabla_{\phi}^2 f_l \cdot v]$$

⇒ Jacobian-vector product (JvP) $\nabla_{\theta\phi}^2 f_l \cdot v$

```

Input: PD sym  $H \in \mathbb{R}^{d \times d}$ , vec  $v \in \mathbb{R}^d$ 
Input: Init  $x_0 \in \mathbb{R}^d$ , prec  $\varepsilon > 0$ , max iters  $n$ 
 $d_0 = r_0 \leftarrow v - Hx_0$ 
for  $i \leftarrow 0, 1, \dots, n$  do
     $\alpha_i \leftarrow (d_i^\top r_i) / (d_i^\top H d_i)$ 
     $x_{i+1} \leftarrow x_i + \alpha_i d_i$ 
     $r_{i+1} \leftarrow r_i - \alpha_i H d_i$ 
     $\beta_{i+1} \leftarrow (r_{i+1}^\top r_{i+1}) / (r_i^\top r_i)$ 
     $d_{i+1} \leftarrow r_{i+1} + \beta_{i+1} d_i$ 
    if  $r_{i+1}^\top r_{i+1} \leq \varepsilon$  then
        return  $x_i$ 
return  $x_{n+1}$ 
    
```


Backward-mode hypergradient $(d\varphi^T/d\theta)^\top v$
for some $v \in \mathbb{R}^{d_l}$

- $\Rightarrow \alpha_T \leftarrow v, g \leftarrow 0 \in \mathbb{R}^{d_u}$
- \Rightarrow for $t = (T - 1) \rightarrow 1$
 - \Rightarrow Compute A_{t+1}, B_{t+1}
 - \Rightarrow Update $g \leftarrow g + B_{t+1} \cdot \alpha_{t+1}$
 - \Rightarrow Update $\alpha_t \leftarrow A_{t+1} \cdot \alpha_{t+1}$
- \Rightarrow Return g

With $\varphi^t \leftarrow \varphi^{t-1} - \beta \nabla_\phi f_l(\theta^k, \varphi^{t-1})$

- $\Rightarrow A_t = \left(I - \beta \nabla_\phi^2 f_l(\theta^k, \varphi^{t-1}) \right),$
- $\Rightarrow B_t = -\beta \nabla_{\theta\phi}^2 f_l(\theta^k, \varphi^{t-1})$
- $\Rightarrow A_t \cdot v = (I \cdot \alpha - \beta \nabla_\phi^2 f_l \cdot v)$
- $\Rightarrow B_t \cdot v = -\beta \nabla_{\theta\phi}^2 f_l(\theta^k, \varphi^{t-1}) \cdot v$

Each iteration needs a HvP and a JvP!

Consider following general operation (Hessian-vector if $\phi = \varphi$):

$$\begin{aligned}\nabla_{\varphi\phi}^2 f(\varphi, \phi) \cdot v &= \left(\nabla_{\varphi} \left(\nabla_{\phi} f(\varphi, \phi) \right) \right)^{\top} \cdot v, \\ f : \mathbb{R}^D \times \mathbb{R}^d &\rightarrow \mathbb{R}, \quad \varphi \in \mathbb{R}^D, \quad \phi \in \mathbb{R}^d, \quad v \in \mathbb{R}^d\end{aligned}\tag{40}$$

Order of operations:

⇒ Compute gradient $\nabla_{\phi} f(\varphi, \phi) \in \mathbb{R}^d$

✗ Compute Jacobian $\nabla_{\varphi} \left(\nabla_{\phi} f(\varphi, \phi) \right)^{\top} \in \mathbb{R}^{D \times d}$

✗ Compute Jacobian-vector product: $\left(\nabla_{\varphi} \left(\nabla_{\phi} f(\varphi, \phi) \right)^{\top} \right) \cdot v \in \mathbb{R}^D$

$$\begin{aligned}\nabla_{\phi\phi}^2 f(\phi, \phi) \cdot v &= \left(\nabla_{\phi} \left(\nabla_{\phi} f(\phi, \phi) \right)^{\top} \right) \cdot v, \\ f: \mathbb{R}^D \times \mathbb{R}^d &\rightarrow \mathbb{R}, \quad \phi \in \mathbb{R}^D, \quad \phi \in \mathbb{R}^d, \quad v \in \mathbb{R}^d\end{aligned}\tag{41}$$

Operations can be re-ordered assuming v is a constant:

- ⇒ Compute gradient $\nabla_{\phi} f(\phi, \phi) \in \mathbb{R}^d$
- ✓ Compute gradient-vector dot-product $\nabla_{\phi} f(\phi, \phi)^{\top} v \in \mathbb{R}$ (scalar)
- ⇒ Compute gradient $\nabla_{\phi} \left(\nabla_{\phi} f(\phi, \phi)^{\top} v \right) \in \mathbb{R}^D$

$$\nabla_{\phi, \phi}^2 f(\phi, \phi) \cdot v = \underbrace{\left(\nabla_{\phi} \underbrace{\left(\nabla_{\phi} f(\phi, \phi) \right)^{\top}}_{1 \times d} \right)}_{D \times d} \cdot \underbrace{v}_{d \times 1} = \nabla_{\phi} \underbrace{\left(\left(\nabla_{\phi} f(\phi, \phi) \right)^{\top} v \right)}_{\text{scalar}},\tag{42}$$

1 Bi-level Problem Variants & Algorithms

- Handling Constraints
- Non-Singleton LL
- Specialized solvers
- Multi-objective Bi-level

2 Implementation Details

3 Further Reading

Method	Single-loop	$\Theta \subset \mathbb{R}^{d_u}$	$\Phi \subset \mathbb{R}^{d_l}$	Hessian-free	Multi-obj	Min-max
BSA [Ghadimi and Wang, 2018]	✗	✓	✗	✗	✗	✗
TTSA [Hong et al., 2020, 2023]	✓	✓	✗	✗	✗	✗
StocBio [Ji et al., 2021]	✗	✗	✗	✗	✗	✗
MRBO [Yang et al., 2021]	✓	✗	✗	✗	✗	✗
VRBO [Yang et al., 2021]	✗	✗	✗	✗	✗	✗
ALSET [Chen et al., 2021]	✓	✗	✗	✗	✗	✗
BSG [Giovannelli et al., 2021]	✗	✓	✓	✗	✗	✗
SignSGD+UL [Fan et al., 2021]	✗	✓	✗	✓	✗	✗
PDBO [Sow et al., 2022a]	✗	✓	✓	✓	✗	✓
STABLE [Chen et al., 2022]	✓	✓	✗	✗	✗	✗
PZOBO [Sow et al., 2022b]	✗	✗	✗	✓	✗	✗
MORBiT [Gu et al., 2022, 2023]	✓	✓	✗	✗	✓	✓
MMB [Hu et al., 2022]	✓	✗	✗	✗	✗	✓

- 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>.
- Mingyi Hong, Hoi-To Wai, Zhaoran Wang, and Zhuoran Yang. A two-timescale stochastic algorithm framework for bilevel optimization: Complexity analysis and application to actor-critic. *SIAM Journal on Optimization*, 33(1):147–180, 2023. URL <https://epubs.siam.org/doi/abs/10.1137/20M1387341>.
- 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. URL <https://proceedings.mlr.press/v151/chen22e/chen22e.pdf>.
- Tommaso Giovannelli, Griffin Kent, and Luis Nunes Vicente. Inexact bilevel stochastic gradient methods for constrained and unconstrained lower-level problems. *arXiv preprint arXiv:2110.00604*, 2021. URL <https://arxiv.org/pdf/2110.00604.pdf>.
- Daouda Sow, Kaiyi Ji, Ziwei Guan, and Yingbin Liang. A primal-dual approach to bilevel optimization with multiple inner minima. *arXiv preprint arXiv:2203.01123*, 2022a. URL <https://arxiv.org/pdf/2203.01123.pdf>.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pages 560–569. PMLR, 2018. URL <http://proceedings.mlr.press/v80/bernstein18a/bernstein18a.pdf>.
- Chen Fan, Parikshit Ram, and Sijia Liu. Sign-MAML: Efficient model-agnostic meta-learning by SignSGD. In *Fifth Workshop on Meta-Learning at the Conference on Neural Information Processing Systems*, 2021. URL <https://arxiv.org/pdf/2109.07497.pdf>.
- Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. Forward and reverse gradient-based hyperparameter optimization. In *International Conference on Machine Learning*, pages 1165–1173. PMLR, 2017. URL <http://proceedings.mlr.press/v70/franceschi17a/franceschi17a.pdf>.
- Alex Gu, Songtao Lu, Parikshit Ram, and Lily Weng. Robust multi-objective bilevel optimization with applications in machine learning. In *INFORMS Annual Meeting*, 2022.
- Alex Gu, Songtao Lu, Parikshit Ram, and Tsui-Wei Weng. Min-max multi-objective bilevel optimization with applications in robust machine learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/pdf?id=PvDY71zKsvP>.
- Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.
- 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.
- Daouda Sow, Kaiyi Ji, and Yingbin Liang. On the convergence theory for hessian-free bilevel algorithms. *Advances in Neural Information Processing Systems*, 35:4136–4149, 2022b. URL <https://openreview.net/pdf?id=suHUIr7dV5n>.
- 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.