

# FORMATTING INSTRUCTIONS FOR ICLR 2023

## CONFERENCE SUBMISSIONS

**Anonymous authors**

Paper under double-blind review

### ABSTRACT

TODO

## 1 INTRODUCTION

TODO

## 2 METHOD

### 2.1 BACKGROUND

We approach the problem of gradient-based hyperparameter optimization as follows. Given a bilevel optimization problem

$$\alpha^* = \arg \min_{\alpha} \mathcal{L}_2(\mathbf{w}^*, \alpha), \quad (1)$$

$$\text{s.t. } \mathbf{w}^* = \arg \min_{\mathbf{w}} \mathcal{L}_1(\mathbf{w}, \alpha). \quad (2)$$

TODO: introduce the functions and the variables.

Suppose that the lower level problem equation 1 is solved by an optimizer which takes the following form

$$\mathbf{w}_{t+1}(\alpha) = \Phi(\mathbf{w}_t, \alpha). \quad (3)$$

Now we derive a gradient of  $\mathcal{L}_2(\mathbf{w}_T(\alpha), \alpha)$  w.r.t.  $\alpha$ . In fact, it could be done by using the chain rule:

$$\underbrace{\frac{d\mathcal{L}_2}{d\alpha}}_{\mathbf{g}_{\text{FO}}} = \frac{\partial \mathcal{L}_2}{\partial \alpha} + \underbrace{\sum_{t=1}^T \mathbf{B}_t \mathbf{A}_{t+1} \dots \mathbf{A}_T \frac{\partial \mathcal{L}_2(\mathbf{w}_T, \alpha)}{\partial \mathbf{w}}}_{\mathbf{g}_{\text{SO}}}, \quad (4)$$

where  $\mathbf{B}_t = \frac{\partial \Phi(\mathbf{w}_{t-1}, \alpha)}{\partial \alpha}$ , and  $\mathbf{A}_t = \frac{\partial \Phi(\mathbf{w}_{t-1}, \alpha)}{\partial \mathbf{w}}$ . However, true hypergradient computation is costly in terms of memory (citation).

### 2.2 HYPERGRADIENT COMPUTATION

In this section we provide an approximation of the true hypergradient and perform analysis of approximation exactness.

Consider the  $t$ -th step of the inner optimization performed by a gradient descent, i.e.  $\Phi(\mathbf{w}_t, \alpha) = \mathbf{w}_t - \eta \frac{\partial \mathcal{L}_1(\mathbf{w}_t, \alpha)}{\partial \mathbf{w}}$ . We motivate the approximation by the fact that we aimed to approximate the  $t$ -th term of  $\mathbf{g}_{\text{SO}}$  equation 4. Since  $\mathbf{A}_{t+1}$  are unknown at the timestamp  $t$ , we approximate the product  $\mathbf{A}_{t+1} \dots \mathbf{A}_T \approx \gamma^{T-t} \mathbf{I}$ , where  $\gamma \in \mathbb{R}_+$ . Additionally, for the same reason we approximate the gradient at the last timestamp  $\frac{\partial \mathcal{L}_2(\mathbf{w}_T, \alpha)}{\partial \mathbf{w}} \approx \frac{\partial \mathcal{L}_2(\mathbf{w}_t, \alpha)}{\partial \mathbf{w}}$ . Finally, the proposed approximation is as follows

$$\hat{\mathbf{g}}_{\text{SO}} = \sum_{t=1}^T \mathbf{B}_t \frac{\partial \mathcal{L}_2(\mathbf{w}_t, \alpha)}{\partial \mathbf{w}} \gamma^{T-t}. \quad (5)$$

Note that in effect the approximation equation 5 is a moving average of greedy hypergradients (cite).

Now we provide analysis of exactness of the provided approximation. Before doing this, we formulate a list of assumptions:

**Assumption 1** 1. Let  $\mathcal{L}_1(\cdot, \alpha)$  and  $\mathcal{L}_2(\cdot, \alpha)$  be  $L$ -smooth and  $\mu$ -strongly convex for any  $\alpha$

2. Let  $\frac{\partial^2 \mathcal{L}_1(\cdot, \alpha)}{\partial \mathbf{w} \partial \mathbf{w}^\top}$  be  $H_w$ -Lipschitz for any  $\alpha$ .

3. Let  $1 - \eta L \leq \gamma \leq 1 - \eta \mu$

4.  $\left\| \frac{\partial \mathcal{L}_1(\mathbf{w}, \alpha)}{\partial \alpha \partial \mathbf{w}^\top} \right\| \leq B$  for any pair  $(\mathbf{w}, \alpha)$ .

**Theorem 1** Let  $\mathcal{L}_1(\cdot, \cdot)$  and  $\mathcal{L}_2(\cdot, \cdot)$  satisfy Assumption 1. Then the following is true

$$\begin{aligned} \|\mathbf{g}_{so} - \hat{\mathbf{g}}_{so}\| &\leq \frac{2LB\|\mathbf{w}_0 - \mathbf{w}_*\|\sqrt{1 - \eta\mu}^T}{\sqrt{1 - \eta\mu}^{-1} - 1} + \\ &B \left\| \frac{\mathcal{L}_2(\mathbf{w}_T, \alpha)}{\partial \mathbf{w}} \right\| \left\{ \frac{1}{\eta} \left( \frac{1}{\mu} - \frac{1}{L} + \frac{1}{L}(1 - \eta\mu)^T \right) + 2\eta H_w \left[ (T-1)\sqrt{1 - \eta\mu}^T - \frac{\sqrt{1 - \eta\mu}^{T-1} - (1 - \eta\mu)^T}{\sqrt{1 - \eta\mu}^{-1} - 1} \right] \right\} \end{aligned}$$

Note that in a general case the approximated gradient is inexact. Notably, when  $L = \mu$  and thus  $H_w = 0$ , we get a simplified upper bound:

$$\|\mathbf{g}_{so} - \hat{\mathbf{g}}_{so}\| \leq \frac{2LB\|\mathbf{w}_0 - \mathbf{w}_*\|\sqrt{1 - \eta\mu}^T}{\sqrt{1 - \eta\mu}^{-1} - 1}. \quad (6)$$

Now we are trying to prove the fact that the bound vanishes if the number of solved inner optimizations tends to infinity. The proof has not been completed yet.

**Theorem 2** In the assumptions above  $\|\mathbf{w}_0^{(k)} - \mathbf{w}_*^k\| \rightarrow_{k \rightarrow \infty} 0$ , where  $k$  is the current epoch. Therefore, the proposed hypergradient is asymptotically exact under the mentioned assumptions.

Furthermore, we provide analysis of a sufficient descent direction. We first formulate an additional list of assumptions.

**Assumption 2** 1. Let  $\mathcal{L}_1(\cdot, \alpha)$  and  $\mathcal{L}_2(\cdot, \alpha)$  be  $M$ -Lipschitz for any  $\alpha$ .

2. Let  $\frac{\partial \mathcal{L}_1(\cdot, \alpha)}{\partial \alpha \partial \mathbf{w}^\top}$  is  $M_b$ -Lipschitz for any  $\alpha$

3. Let  $(\frac{\partial \mathcal{L}_1(\cdot, \alpha)}{\partial \alpha \partial \mathbf{w}^\top})^\top \frac{\partial \mathcal{L}_1(\cdot, \alpha)}{\partial \alpha \partial \mathbf{w}^\top} \succeq \kappa \mathbf{I}$  for any  $(\mathbf{w}, \alpha)$

**Lemma 1** Given an gradient descent update of the outer optimization with a learning rate  $\eta_{out}$ . Suppose that we start the whole optimization from  $(\mathbf{w}_0, \alpha_0)$ . Suppose that Assumption 1 and the first point of Assumption 2 hold. Let also the following conditions are met:

1. There exist  $\delta > 0$  such that  $\mathcal{L}_2(\mathbf{w}_*, \alpha) - \mathcal{L}_2^* - M\|\mathbf{w}_0 - \mathbf{w}_*\|\sqrt{1 - \eta\mu}^T \geq \delta$  for any  $\alpha$

2.  $\sqrt{1 - \eta\mu}^T \leq 1/2$

3.  $\eta_{out} \leq \frac{\|\mathbf{w}_0 - \mathbf{w}_*\|}{2} \left( \frac{LB^2}{(1 - \eta L)^2} \right)^{-1}$

Then the following is true for any  $\alpha$

$$\left\| \frac{\partial \mathcal{L}_2(\mathbf{w}_T, \alpha)}{\partial \mathbf{w}} \right\| \geq \sqrt{\mu\delta}. \quad (7)$$

**Theorem 3** Let Assumptions 1, 2 are satisfied. Additionally, let there exist large enough  $\delta > 0$ . Then there exist  $c > 0$  such that

$$\mathbf{g}_{so}^\top \hat{\mathbf{g}}_{so} \geq c \left\| \frac{\partial \mathcal{L}_2(\mathbf{w}_T, \alpha)}{\partial \mathbf{w}} \right\|^2. \quad (8)$$

In other words, if the outer function does not depend on the hyperparameters  $\mathcal{L}_2 = \mathcal{L}_2(\mathbf{w})$ , then the proposed approximation  $\mathbf{g}_{SO}$  is sufficient descent condition.

I was inspired by "Truncated Back-propagation for Bilevel Optimization" when formulating and proving this theorem.

### 3 EXPERIMENTS

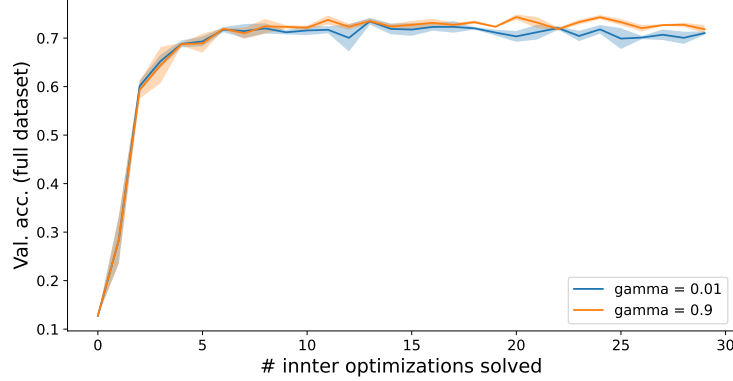


Figure 1: Experimental results on data hypercleaning.

Method	Valid. Acc.	#JVPs
Truncated backpropagation (Lukethina)	72.5	1(1)
DrMAD	69.8	99(2T - 1)
IFT(9, 5)	70.3	50((N + 1)K)
IFT(4, 10)	70.7	50((N + 1)K)
Proposed ( $\gamma = 0.99$ )	<b>73.5*</b>	50(T)

Table 1: Experimental results for data hypercleaning. \* indicates that the results are statistically significant ( $p < 0.05$ ). Note IFT( $N, K$ ) denotes that we performed  $K$ -step online optimization with  $N$  first terms of Neuman series.

### REFERENCES

### A APPENDIX

You may include other additional sections here.