

Stochastic Bilevel Optimization

January 4, 2023

1 Preliminaries

Let us list some useful definitions. We have

$$\nabla \mathcal{L}(\lambda_k) = \nabla_\lambda C(\hat{x}(\lambda_k), \lambda_k) - \nabla_{x\lambda}^2 F(\hat{x}(\lambda_k), \lambda_k)^\top [\nabla_{xx} F(\hat{x}(\lambda_k), \lambda_k)]^{-1} \nabla_x C(\hat{x}(\lambda_k), \lambda_k)$$

$$\tilde{\nabla} \mathcal{L}(\lambda_k) = \nabla_\lambda C(x_k, \lambda_k) - \nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top [\nabla_{xx} F(x_k, \lambda_k)]^{-1} \nabla_x C(x_k, \lambda_k)$$

$$\hat{\nabla} \mathcal{L}(\lambda_k) = \nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q$$

Definition 1.1 (Jensen's Inequality). *Theorem 1 (Jensen's Inequality) Let φ be a convex function on \mathbb{R} and let $X \in L_1$ be integrable. Then*

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$$

2 Problem

We now have the bilevel problem

$$\min_{\lambda \in \mathcal{D} \subset \mathbb{R}^r} \mathcal{L}(\lambda) \triangleq C(\hat{x}(\lambda), \lambda) \tag{1}$$

$$\text{s.t. } \hat{x}(\lambda) = \arg \min_{x \in \mathbb{R}^n} F(x, \lambda). \tag{2}$$

We will denote the sampled terms as follows:

$$\mathbb{E}_\zeta[\nabla_{xx}^2 \tilde{F}(x_k, \lambda_k; \zeta)] = \nabla_{xx}^2 F(x_k, \lambda_k)$$

$$\mathbb{E}_\xi[\nabla_x \tilde{C}(x_k, \lambda_k, \xi)] = \nabla_x \tilde{C}(x_k, \lambda_k)$$

$$\mathbb{E}_\theta[\nabla_\lambda \tilde{C}(x_k, \lambda_k, \theta)] = \nabla_\lambda \tilde{C}(x_k, \lambda_k)$$

$$\mathbb{E}_\kappa[\nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \kappa)] = \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k)$$

3 Algorithm

Algorithm 1 Stochastic HOAG

1: At iteration $k = 1, 2, \dots$, given random samples ξ_i, ζ_j , perform the following:

1. Solve the inner optimization problem up to tolerance ε_k . That is, find x_k such that

$$\mathbb{E} [\|\hat{x}(\lambda_k) - x_k\|] \leq \varepsilon_k$$

2.

$$v_Q = \eta \sum_{q=-1}^{Q-1} \prod_{j=Q-q}^Q \left(I - \eta \nabla_{xx}^2 \tilde{F}(x_k, \lambda_k; \zeta_j) \right) \nabla_x \tilde{C}(x_k, \lambda_k, \xi_0), \quad (3)$$

3. Compute approximate stochastic gradient $\hat{\nabla} \mathcal{L}(\lambda_k)$ as

$$\hat{\nabla} \mathcal{L}(\lambda_k) = \nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q$$

4. Update hyperparameters:

$$\lambda_{k+1} = P_{\mathcal{D}} \left(\lambda_k - \frac{1}{L} \hat{\nabla} \mathcal{L}(\lambda_k) \right).$$

4 Stochastic HOAG

In the following section we adapt the convergence proof of HOAG to the case when all terms are sampled using a single sample.

Assumption 4.1 (Convexity). *The lower-level function $F(x, \lambda)$ is μ strongly-convex w.r.t. x and the total objective function $\mathcal{L}(\lambda) = C(\lambda, \hat{x}(\lambda))$ is nonconvex w.r.t. λ . For the stochastic setting, the same assumptions hold for $F(x, \lambda; \zeta)$ and $\mathcal{L}(\lambda, \zeta)$, respectively.*

Assumption 4.2 (Smoothness). *Let $z = (x, \lambda) \in \mathbb{R}^n \times \mathcal{D}$. The loss function $C(z)$ and $F(z)$ satisfy - The function $C(z)$ is M -Lipschitz, i.e., for any z, z' ,*

$$|C(z) - C(z')| \leq M \|z - z'\|.$$

- $\nabla C(z)$ and $\nabla F(z)$ are L -Lipschitz, i.e., for any z, z' ,

$$\|\nabla C(z) - \nabla C(z')\| \leq L \|z - z'\|,$$

$$\|\nabla F(z) - \nabla F(z')\| \leq L \|z - z'\|.$$

For the stochastic case, the same assumptions hold for $F(z; \xi)$ and $G(z; \zeta)$ for any given ξ and ζ .

Assumption 4.3 (Partial Lipschitz Smoothness). *Let $z = (x, \lambda) \in \mathbb{R}^n \times \mathcal{D}$. Suppose the derivatives $\nabla_{x\lambda} F(z)$ and $\nabla_x^2 F(z)$ are τ - and ρ -Lipschitz, i.e., - For any z, z' , $\|\nabla_{x\lambda} F(z) - \nabla_{x\lambda} F(z')\| \leq \tau \|z - z'\|$. - For any z, z' , $\|\nabla_x^2 F(z) - \nabla_x^2 F(z')\| \leq \rho \|z - z'\|$. For the stochastic case, the same assumptions hold for $\nabla_{x\lambda} F(z; \zeta)$ and $\nabla_x^2 F(z; \zeta)$ for any ζ .*

4.1 Convergence Proof Attempt

Theorem 4.4 (Bounded Gradient Error (Lemma 7 of [2])). *Suppose Assumptions 4.1, 4.2 and 4.3 hold. Then, conditioning on x_k^D and λ_k , we have*

$$\left\| \mathbb{E} \hat{\nabla} \mathcal{L}(\lambda_k) - \nabla \mathcal{L}(\lambda_k) \right\|^2 \leq 2 \left(L + \frac{L^2}{\mu} + \frac{M\tau}{\mu} + \frac{LM\rho}{\mu^2} \right)^2 \|x_k^D - \hat{x}(\lambda_k)\|^2 + \frac{2L^2 M^2 (1 - \eta\mu)^{2Q}}{\mu^2},$$

that is,

$$\left\| \mathbb{E} \hat{\nabla} \mathcal{L}(\lambda_k) - \nabla \mathcal{L}(\lambda_k) \right\| = \mathcal{O}(\epsilon_k), \quad \|x_k^D - \hat{x}(\lambda_k)\| \leq \epsilon_k.$$

Theorem 4.5 (Bounded Variance of $\widehat{\nabla}\mathcal{L}$ (Lemma 8 of [2])). *Suppose Assumptions 4.1, 4.2 and 4.3 hold. Assume all sample sizes are 1. Then, we have*

$$\begin{aligned} \mathbb{E} \left\| \widehat{\nabla}\mathcal{L}(\lambda_k) - \nabla\mathcal{L}(\lambda_k) \right\|^2 &\leq \frac{4L^2M^2}{\mu^2} + \left(\frac{8L^2}{\mu^2} + 2 \right) \frac{M^2}{1} + \frac{16\eta^2L^4M^2}{\mu^2} \frac{1}{B} + \frac{16L^2M^2(1-\eta\mu)^{2Q}}{\mu^2} \\ &\quad + \left(L + \frac{L^2}{\mu} + \frac{M\tau}{\mu} + \frac{LM\rho}{\mu^2} \right)^2 \mathbb{E} \|x_k^D - \hat{x}(\lambda_k)\|^2. \end{aligned}$$

That is,

$$\mathbb{E} \left[\left\| \widehat{\nabla}\mathcal{L}(\lambda_k) - \nabla\mathcal{L}(\lambda_k) \right\|^2 \right] = \mathcal{O}(\mathbb{E} [\|x_k^D - \hat{x}(\lambda_k)\|^2]) \stackrel{?}{=} \mathcal{O}(\epsilon^2).$$

Theorem 4.6 (Global Convergence (Simplified) (SGD step)). *In Algorithm 3, assume that the update*

$$\lambda_{k+1} = P_{\mathcal{D}} \left(\lambda_k - \frac{1}{L} \widehat{\nabla}\mathcal{L}(\lambda_k) \right),$$

is replaced by SGD with approximate gradient:

$$\lambda_{k+1} = \lambda_k - \frac{1}{L} \widehat{\nabla}\mathcal{L}(\lambda_k).$$

Assume also, that $\lambda_k \in \mathcal{D}$ for all $k > 0$. If the sequence ϵ_k obeys

$$\sum_{i=1}^{\infty} \epsilon_i < \infty, \quad \epsilon_k > 0 \quad \forall k \geq 0,$$

Proof. An equivalent condition to $\mathcal{L}(\lambda)$ having Lipschitz continuous gradient is that for any $\alpha, \beta \in \mathcal{D}$:

$$\mathcal{L}(\beta) \leq \mathcal{L}(\alpha) + \nabla\mathcal{L}(\alpha)^\top(\beta - \alpha) + \frac{L}{2} \|\beta - \alpha\|^2. \quad (4)$$

Substituting for $\alpha = \lambda_k, \beta = \lambda_{k+1} = \lambda_k - \frac{1}{L} \widehat{\nabla}\mathcal{L}(\lambda_k)$, we add and subtract $\widehat{\nabla}\mathcal{L}(\lambda_k)^\top(\lambda_k - \lambda_{k+1})$ to get

$$\mathcal{L}(\lambda_{k+1}) \leq \mathcal{L}(\lambda_k) - (\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k))^\top(\lambda_k - \lambda_{k+1}) - \widehat{\nabla}\mathcal{L}(\lambda_k)^\top(\lambda_k - \lambda_{k+1}) + \frac{L}{2} \|\lambda_k - \lambda_{k+1}\|^2 \quad (5)$$

$$= \mathcal{L}(\lambda_k) - (\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k))^\top(\lambda_k - \lambda_{k+1}) - \frac{1}{L} \|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2 + \frac{1}{2L} \|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2 \quad (\text{using } \lambda_k - \lambda_{k+1} = \frac{1}{L} \widehat{\nabla}\mathcal{L}(\lambda_k)) \quad (6)$$

$$= \mathcal{L}(\lambda_k) - (\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k))^\top(\lambda_k - \lambda_{k+1}) - \frac{1}{2L} \|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2. \quad (7)$$

Finally, by Cauchy-Schwarz, we have

$$\mathcal{L}(\lambda_{k+1}) \leq \mathcal{L}(\lambda_k) + \|\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k)\| \|\lambda_k - \lambda_{k+1}\| - \frac{1}{2L} \|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2. \quad (8)$$

Taking expectations of both sides, and conditioning on λ_k , we get

$$\begin{aligned} \mathbb{E}\mathcal{L}(\lambda_{k+1}) &\leq \mathcal{L}(\lambda_k) + \mathbb{E} \left[\|\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k)\| \|\lambda_k - \lambda_{k+1}\| \right] - \frac{1}{2L} \mathbb{E} [\|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2] \\ &= \mathcal{L}(\lambda_k) + \mathbb{E} \left[\|\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k)\| \left\| \frac{1}{L} \widehat{\nabla}\mathcal{L}(\lambda_k) \right\| \right] - \frac{1}{2L} \mathbb{E} [\|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2] \\ &\leq \mathcal{L}(\lambda_k) + \frac{1}{L} \sqrt{\mathbb{E} [\|\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k)\|^2]} \sqrt{\mathbb{E} [\|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2]} - \frac{1}{2L} \mathbb{E} [\|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2], \end{aligned} \quad (9)$$

where the last line is given by the Cauchy-Schwarz inequality.

4.2 Approach I

That root term looks nasty to work with - let's get rid of it. Denote

$$X = \mathbb{E} \left[\|\nabla \mathcal{L}(\lambda_k) - \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right], \quad Y = \mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right].$$

Then we have, that

$$\begin{aligned} & \frac{1}{L} \sqrt{\mathbb{E} \left[\|\nabla \mathcal{L}(\lambda_k) - \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right]} \sqrt{\mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right]} - \frac{1}{2L} \mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right] = \frac{1}{L} \sqrt{XY} - \frac{1}{2L} Y \\ & = \frac{1}{L} X - \frac{1}{4L} Y - \frac{1}{L} \left(\sqrt{X} - \frac{1}{2} \sqrt{Y} \right)^2. \end{aligned}$$

Thus, going back to equation (9), we get

$$\begin{aligned} \mathbb{E} \mathcal{L}(\lambda_{k+1}) & \leq \mathcal{L}(\lambda_k) + \frac{1}{L} \sqrt{\mathbb{E} \left[\|\nabla \mathcal{L}(\lambda_k) - \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right]} \sqrt{\mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right]} - \frac{1}{2L} \mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right] \\ & = \mathcal{L}(\lambda_k) + \frac{1}{L} \mathbb{E} \left[\|\nabla \mathcal{L}(\lambda_k) - \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right] - \frac{1}{4L} \mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right] - \frac{1}{L} \left(\sqrt{\mathbb{E} \left[\|\nabla \mathcal{L}(\lambda_k) - \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right]} - \frac{1}{2} \sqrt{\mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right]} \right)^2 \\ & \leq \mathcal{L}(\lambda_k) + \frac{1}{L} \mathbb{E} \left[\|\nabla \mathcal{L}(\lambda_k) - \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right] - \frac{1}{4L} \mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right]. \\ & \leq \mathcal{L}(\lambda_k) + \frac{M}{L} \mathbb{E} \left[\|x_k^D - \hat{x}(\lambda_k)\|^2 \right] - \frac{1}{4L} \mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right]. \end{aligned}$$

where the last line follows from theorem 4.5. Rewriting, we get

$$\frac{1}{4L} \mathbb{E} \left[\|\widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \right] \leq \mathcal{L}(\lambda_k) - \mathbb{E} \mathcal{L}(\lambda_{k+1}) + \frac{M}{L} \mathbb{E} \left[\|x_k^D - \hat{x}(\lambda_k)\|^2 \right] \quad (10)$$

4.3 Approach II

An equivalent condition to $\mathcal{L}(\lambda)$ having Lipschitz continuous gradient is that for any $\alpha, \beta \in \mathcal{D}$:

$$\mathcal{L}(\beta) \leq \mathcal{L}(\alpha) + \nabla \mathcal{L}(\alpha)^\top (\beta - \alpha) + \frac{L}{2} \|\beta - \alpha\|^2. \quad (11)$$

Substituting for $\alpha = \lambda_k, \beta = \lambda_{k+1} = \lambda_k - \frac{1}{L} \widehat{\nabla} \mathcal{L}(\lambda_k)$,

$$\mathcal{L}(\lambda_{k+1}) \leq \mathcal{L}(\lambda_k) + \nabla \mathcal{L}(\lambda_k)^\top \left(-\frac{1}{L} \widehat{\nabla} \mathcal{L}(\lambda_k) \right) + \frac{L}{2} \left\| -\frac{1}{L} \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|^2.$$

Taking expectation, conditioning on λ_k ,

$$\begin{aligned} \mathbb{E} [\mathcal{L}(\lambda_{k+1})] & \leq \mathcal{L}(\lambda_k) - \frac{1}{L} \nabla \mathcal{L}(\lambda_k)^\top \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] + \frac{1}{2L} \mathbb{E} \left[\left\| \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \right]. \\ & = \mathcal{L}(\lambda_k) - \frac{1}{L} \left(\nabla \mathcal{L}(\lambda_k) - \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right)^\top \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] - \frac{1}{L} \left\| \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\|^2 + \frac{1}{2L} \mathbb{E} \left[\left\| \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \right] \\ & \leq \mathcal{L}(\lambda_k) + \frac{1}{L} \left\| \nabla \mathcal{L}(\lambda_k) - \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\| \left\| \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\| - \frac{1}{L} \left\| \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\|^2 + \frac{1}{2L} \mathbb{E} \left[\left\| \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \right] \\ & \leq \mathcal{L}(\lambda_k) + \frac{1}{L} \left\| \nabla \mathcal{L}(\lambda_k) - \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\| \left\| \mathbb{E} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\| + \frac{1}{2L} \text{Var}(\widehat{\nabla} \mathcal{L}(\lambda_k)) - \frac{1}{2L} \left\| \mathbb{E} \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \end{aligned}$$

Now, we are interested in controlling $\text{Var}(\widehat{\nabla}\mathcal{L}(\lambda_k))$ through $\text{Var}(v_Q)$, where

$$\begin{aligned}
v_Q &= \eta \sum_{q=-1}^{Q-1} \prod_{j=Q-q}^Q \left(I - \eta \nabla_{xx}^2 \tilde{F}(x_k, \lambda_k; \mathcal{B}_j) \right) \nabla_x \tilde{C}(x_k, \lambda_k, \mathcal{S}_C), \\
\text{Var}(v_Q) &= \mathbb{E} \|v_Q - \mathbb{E}v_Q\|^2 \\
\text{Var}(\widehat{\nabla}\mathcal{L}(\lambda_k)) &= \mathbb{E} \|\widehat{\nabla}\mathcal{L}(\lambda_k) - \mathbb{E}\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2 = \\
&= \mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q - (\nabla_\lambda C(x_k, \lambda_k) - \nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E}v_Q)\|^2 \\
&= \mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k)\|^2 + \mathbb{E} \|\nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E}v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q\|^2 - \\
&\quad - 2\mathbb{E} \left[\underbrace{\left(\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k) \right)}_{\mathbf{0} \text{ in expectation}}^\top \left(\nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E}v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q \right) \right] \\
&= \mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k)\|^2 + \mathbb{E} \|\nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E}v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q\|^2 \\
&= \text{Var} \left(\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) \right) + \text{Var} \left(\nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q \right)
\end{aligned} \tag{12}$$

4.4 Rest of original proof

Now, by Theorem 4.4, $\|\nabla\mathcal{L}(\lambda_k) - \mathbb{E}\widehat{\nabla}\mathcal{L}(\lambda_k)\| = \mathcal{O}(\epsilon_k)$. As \mathcal{D} is bounded (by Heine-Borel), we have that $\exists M > 0$, such that

$$\|\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k)\| \|\lambda_k - \lambda_{k+1}\| < M\epsilon_k.$$

Applying this to (8), we get

$$\mathcal{L}(\lambda_{k+1}) \leq \mathcal{L}(\lambda_k) + M\epsilon_k - \frac{1}{2L} \|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2. \tag{13}$$

Rewriting, we get

$$\frac{1}{2L} \|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2 \leq \mathcal{L}(\lambda_k) - \mathcal{L}(\lambda_{k+1}) + M\epsilon_k. \tag{14}$$

By the extreme-value theorem, since L is defined on a compact set and has continuous derivatives, it has a lower bound K . Thus, taking the sum of (14) for $k = m$ to ∞ , we get

$$\frac{1}{2L} \sum_{k=m}^{\infty} \|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2 \leq \mathcal{L}(\lambda_m) - K + M \sum_{k=m}^{\infty} \epsilon_k. \tag{15}$$

$\{\epsilon_k\}_{k=1}^{\infty}$ is summable by assumption, thus $\sum_{k=m}^{\infty} \epsilon_k < \infty$, $\epsilon_k \rightarrow 0$ and the RHS of (15) is finite. Hence the LHS of (15) must also be finite. Hence $\|\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2 \rightarrow 0 \iff \|\widehat{\nabla}\mathcal{L}(\lambda_k)\| \rightarrow 0$. Recall, that $\|\nabla\mathcal{L}(\lambda_k) - \widehat{\nabla}\mathcal{L}(\lambda_k)\| = \mathcal{O}(\epsilon_k)$, hence

$$\|\nabla\mathcal{L}(\lambda_k)\| \leq \|\widehat{\nabla}\mathcal{L}(\lambda_k)\| + \mathcal{O}(\epsilon_k) \xrightarrow{k \rightarrow \infty} 0.$$

□

References

- [1] Bottou, L., Curtis, F., & Nocedal, J.. (2016). Optimization Methods for Large-Scale Machine Learning.
- [2] Ji, K. Yang, J. & Liang, Y. Bilevel optimization: Nonasymptotic analysis and faster algorithms. *International Conference on Machine Learning (ICML)*, 2021.
- [3] Pedregosa, F. Hyperparameter optimization with approximate gradient. *Proceedings of The 33rd International Conference on Machine Learning*, PMLR 48:737-746, 2016. Available from <https://proceedings.mlr.press/v48/pedregosa16.html>.