

# Stochastic Bilevel Optimization

January 31, 2023

## 1 Problem

We now have the bilevel problem

$$\begin{aligned} \min_{\lambda \in \mathcal{D} \subset \mathbb{R}^r} \mathcal{L}(\lambda) &\triangleq C(\hat{x}(\lambda), \lambda) \\ \text{s.t. } \hat{x}(\lambda) &= \arg \min_{x \in \mathbb{R}^n} F(x, \lambda). \end{aligned} \tag{1}$$

We will denote the sampled terms as follows:

$$\begin{aligned} \mathbb{E}_{\zeta}[\tilde{F}(x_k, \lambda_k; \zeta)] &= F(x_k, \lambda_k) \\ \mathbb{E}_{\xi}[\tilde{C}(x_k, \lambda_k, \xi)] &= C(x_k, \lambda_k) \end{aligned}$$

## 2 Algorithm

---

### Algorithm 1 Stochastic HOAG

---

1: At iteration  $k = 1, 2, \dots$ , given random samples  $\xi_i, \zeta_j$ , stepsize  $\nu_k$ , perform the following:

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

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

2.

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

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:

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


---

## 3 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 3.1** (Convexity). *The lower-level function  $F(x, \lambda)$  is  $\mu$  strongly-convex w.r.t.  $x$  and the total objective function  $\mathcal{L}(\lambda) = C(\lambda, \hat{x}(\lambda))$  is nonconvex w.r.t.  $\lambda$ . For the stochastic setting, the same assumptions hold for  $F(x, \lambda; \zeta)$  and  $\mathcal{L}(\lambda, \zeta)$ , respectively.*

**Assumption 3.2** (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  $\xi$  and  $\zeta$ .

**Assumption 3.3** (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  $\tau$  - and  $\rho$  - 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  $\zeta$ .

**Assumption 3.4** (Bounded Gradient). Assume that the partial gradient  $\nabla_{x\lambda}^2 F$  is bounded in norm, i.e.  $\|\nabla_{x\lambda}^2 F\| \leq K$ .

**Assumption 3.5** (Lower bound on objective.). The sequence of iterates  $\{\lambda_k\}$  is contained in an open set over which  $\mathcal{L}$  is bounded below by a scalar  $\mathcal{L}_{\inf}$ .

### 3.1 Preliminary Results

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

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

**Lemma 3.7** (Neumann Series). For non-singular  $A \in \mathbb{R}^{n \times n}$ ,

$$A^{-1} = \sum_{i=0}^{\infty} (I - A)^i, \quad A \succ 0, \|A\| < 1. \quad (4)$$

Firstly, we will present an immediate consequence of the assumptions in the previous section.

**Proposition 3.8** (Bounded variance of  $\nabla \tilde{C}, \nabla \tilde{F}, \nabla_{x\lambda}^2 \tilde{F}, \nabla_{xx}^2 \tilde{F}$ . Lemma 1 in [2]). Suppose, Assumption 3.2 holds. Then for any  $z = (x, \lambda), \zeta$ ,

$$\mathbb{E}_\zeta \|\nabla \tilde{C}(z, \zeta) - \nabla C(z)\|^2 \leq M^2$$

$$\mathbb{E}_\zeta \|\nabla_{x\lambda}^2 \tilde{F}(z, \zeta) - \nabla_{x\lambda}^2 F(z)\|^2 \leq L^2$$

$$\mathbb{E}_\zeta \|\nabla_{xx}^2 \tilde{F}(z, \zeta) - \nabla_{xx}^2 F(z)\|^2 \leq L^2$$

Note that

$$v_Q = \eta \sum_{q=-1}^{Q-1} \prod_{j=Q-q}^Q \left( I - \eta \nabla_{xx}^2 \tilde{F}(x_k, \lambda_k; \zeta_j) \right) \nabla_x \tilde{C}(x_k, \lambda_k, \xi_0), \quad (5)$$

where we assume  $\prod_{j=Q+1}^Q (\cdot) = I$ . From this we easily get

$$\mathbb{E}[v_Q] = \eta \sum_{i=0}^Q [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \nabla_x C(x_k, \lambda_k).$$

Denote by  $\mathbb{E}[v_\infty]$ :

$$\mathbb{E}[v_\infty] = [\nabla_{xx}^2 F(x_k, \lambda_k)]^{-1} \nabla_x C(x_k, \lambda_k)$$

**Proposition 3.9** (Bound on  $\|\mathbb{E}v_Q\|$ ). *Suppose Assumptions 3.1, 3.2 hold. Then*

$$\|\mathbb{E}v_Q\| \leq \frac{M}{\mu}(1 - (1 - \eta\mu)^{Q+1})$$

**Proposition 3.10** (Bound on  $\text{Var}(v_Q)$ ). *Suppose Assumptions 3.1, 3.2 and 3.4 hold. Denote the condition number as  $\kappa = \frac{L}{\mu}$ . Choose  $\eta$ , such that  $\eta\mu < 1$ . Then we have that*

$$\text{Var}(v_Q) = \mathbb{E}\|v_Q - \mathbb{E}(v_Q)\|^2 \leq 2\eta^3 M^2 L\kappa + \frac{2\eta M^2}{\mu}.$$

**Proposition 3.11** (Bound on  $\text{Var}(\widehat{\nabla}\mathcal{L})$ ). *Suppose Assumptions 3.1, 3.2 and 3.4 hold. Choose  $\eta$ , such that  $\eta\mu < 1$ . Denote the condition number as  $\kappa = \frac{L}{\mu}$ . Then the variance of the approximate hypergradient satisfies*

$$\begin{aligned} \text{Var}(\widehat{\nabla}\mathcal{L}(\lambda_k)) &= \mathbb{E}\|\widehat{\nabla}\mathcal{L}(\lambda_k) - \mathbb{E}\widehat{\nabla}\mathcal{L}(\lambda_k)\|^2 \leq \\ &\text{Var}(v_Q)(L^2 + K^2) + L^2\|\mathbb{E}v_Q\|^2 + M^2 \leq \\ &\left(2\eta^3 M^2 L\kappa + \frac{2\eta M^2}{\mu}\right)(L^2 + K^2) + \kappa^2 M^2 + M^2. \end{aligned}$$

In particular, we have that if  $\text{Var}(v_Q)$  is bounded, then so is  $\text{Var}(\widehat{\nabla}\mathcal{L}(\lambda_k))$ .

**Proposition 3.12** (Bound on  $\|\mathbb{E}_{\lambda_k}\widehat{\nabla}\mathcal{L}(\lambda_k) - \nabla\mathcal{L}(\lambda_k)\|$ . Lemma 7 of [2]). *Let*

$$T_4 = \sqrt{2}\left(L + \frac{L^2}{\mu} + \frac{M\tau}{\mu} + \frac{LM\rho}{\mu^2}\right), \quad T_5 = \sqrt{2}\frac{LM(1 - \eta\mu)^Q}{\mu}. \quad (6)$$

*Then we have that*

$$\|\mathbb{E}_{\lambda_k}\widehat{\nabla}\mathcal{L}(\lambda_k) - \nabla\mathcal{L}(\lambda_k)\|_2 \leq T_4\epsilon_k + T_5.$$

### 3.2 Convergence Proof

**Theorem 3.13** (Global Convergence (SGD step)). *In Algorithm 2, assume that the stepsize  $\nu_k$  is chosen such that*

$$\sum_{k=1}^{\infty} \nu_k = \infty, \quad \sum_{k=1}^{\infty} \nu_k^2 < \infty.$$

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

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

*then we have*

$$\liminf_{k \rightarrow \infty} \mathbb{E} \left[ \|\nabla\mathcal{L}(\lambda_k)\|_2^2 \right] = 0.$$

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

$$\mathcal{L}(\beta) \leq \mathcal{L}(\alpha) + \nabla\mathcal{L}(\alpha)^\top(\beta - \alpha) + \frac{L}{2}\|\beta - \alpha\|^2. \quad (7)$$

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

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

Taking expectation, conditioning on  $\lambda_k$ ,

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

$$\begin{aligned}
&= \mathcal{L}(\lambda_k) - \nu_k \left( \nabla \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right)^\top \mathbb{E}_{\lambda_k} [\widehat{\nabla} \mathcal{L}(\lambda_k)] - \nu_k \left\| \mathbb{E}_{\lambda_k} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\|^2 + \frac{L\nu_k^2}{2} \mathbb{E}_{\lambda_k} \left[ \left\| \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \right] \\
&\leq \mathcal{L}(\lambda_k) + \nu_k \left\| \nabla \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\| \left\| \mathbb{E}_{\lambda_k} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\| - \nu_k \left\| \mathbb{E}_{\lambda_k} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\|^2 + \frac{L\nu_k^2}{2} \mathbb{E}_{\lambda_k} \left[ \left\| \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \right] \\
&= \mathcal{L}(\lambda_k) + \nu_k \left\| \nabla \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\| \left\| \mathbb{E}_{\lambda_k} [\widehat{\nabla} \mathcal{L}(\lambda_k)] \right\| + \frac{L\nu_k^2}{2} \text{Var}(\widehat{\nabla} \mathcal{L}(\lambda_k)) - \left( \nu_k - \frac{L\nu_k^2}{2} \right) \left\| \mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|^2,
\end{aligned} \tag{8}$$

Let  $T_{45} = T_4\epsilon_k + T_5$ , where  $T_4, T_5$  are as in (6) of Proposition 3.12. From Proposition 3.12 and by assumption on  $\epsilon_k$ , we have that

$$\left\| \mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k) - \nabla \mathcal{L}(\lambda_k) \right\|_2 \leq T_{45}. \tag{9}$$

By reverse triangle inequality, we have that

$$\left| \left\| \mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|_2 - \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2 \right| \leq T_{45}.$$

Hence,

$$\left\| \mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|_2 \leq T_{45} + \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2$$

Note, that from this we have

$$\left\| \mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|_2 \leq T_{45} + 1 + \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2 \tag{10}$$

Furthermore,

$$\begin{aligned}
\left\| \mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k) \right\|_2^2 &\leq T_{45}^2 + 2T_{45} \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2 + \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2 \leq T_{45}^2 + 2T_{45} (1 + \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2) + \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2 \\
&= T_{45}^2 + 2T_{45} + (2T_{45} + 1) \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2.
\end{aligned} \tag{11}$$

Finally, by proposition 3.11, we have that

$$\text{Var} \left( \widehat{\nabla} \mathcal{L}(\lambda_k) \right) \leq T_2 = \left( 2\eta^3 M^2 L \kappa + \frac{2\eta M^2}{\mu} \right) (L^2 + K^2) + \kappa^2 M^2 + M^2. \tag{12}$$

Set  $M_1 = T_{45}^2 + 2T_{45}$ ,  $M_2 = 2T_{45} + 1$ . Substituting (9),(10), (11) and (12) into (8), and rearranging, we get

$$\begin{aligned}
&\left( \nu_k - \frac{L\nu_k^2}{2} \right) (M_1 + M_2 \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2) - \nu_k ((T_4\epsilon + T_5) \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2 + T_5^2 + T_5) \leq \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\mathcal{L}(\lambda_{k+1})] + \\
&\quad + \nu_k \epsilon_k \left( \underbrace{T_4^2 T_\epsilon + 2T_4 T_5 + T_4}_{\triangleq T_6} \right) + \frac{T_2 \nu_k^2}{2L}
\end{aligned} \tag{13}$$

Now, let us inspect the left-hand-side of (13). Expanding  $M_1, M_2$  and collecting like terms, we get

$$\begin{aligned}
&\left( \nu_k - \frac{L\nu_k^2}{2} \right) (M_1 + M_2 \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2) - \nu_k ((T_4\epsilon + T_5) \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2 + T_5^2 + T_5) = \\
&= \left( \underbrace{\nu_k (2T_{45} + 1 - T_5) - \frac{M_2 L \nu_k^2}{2} - T_4 \nu_k \epsilon_k}_{\triangleq \gamma_k} \right) \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2 + \nu_k \epsilon_k (T_4\epsilon + 2T_4 T_5 + 2T_4) + \nu_k T_5 \geq \\
&\geq \gamma_k \left\| \nabla \mathcal{L}(\lambda_k) \right\|_2^2 + \nu_k \epsilon_k (T_4\epsilon + 2T_4 T_5 + 2T_4).
\end{aligned}$$

We can now rewrite the inequality (13) as

$$\gamma_k \|\nabla \mathcal{L}(\lambda_k)\|_2^2 \leq \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\mathcal{L}(\lambda_{k+1})] + \frac{T_2 \nu_k^2}{2L} + \underbrace{\nu_k \epsilon_k (T_6 - T_4 \epsilon - 2T_4 T_5 - 2T_4)}_{\triangleq T_7} \quad (14)$$

Note that, by Cauchy-Schwarz,  $\sum_{k=1}^{\infty} \nu_k \epsilon_k \leq \sqrt{\sum_{k=1}^{\infty} \nu_k^2} \sqrt{\sum_{k=1}^{\infty} \epsilon_k^2} < \infty$ . As such,  $\gamma_k$  is summable. Summing (14) for  $k = 1$  to  $\infty$ , we get

$$\sum_{k=1}^{\infty} \gamma_k \|\nabla \mathcal{L}(\lambda_k)\|_2^2 \leq \sum_{k=1}^{\infty} (\mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\mathcal{L}(\lambda_{k+1})]) + \frac{T_2}{2L} \sum_{k=1}^{\infty} \nu_k^2 + T_7 \sum_{k=1}^{\infty} \nu_k \epsilon_k.$$

Taking total expectation and telescoping, we get

$$\sum_{k=1}^{\infty} \gamma_k \mathbb{E} \|\nabla \mathcal{L}(\lambda_k)\|_2^2 \leq \mathcal{L}(\lambda_1) - \mathcal{L}_{inf} + \frac{T_2}{2L} \sum_{k=1}^{\infty} \nu_k^2 + T_7 \sum_{k=1}^{\infty} \nu_k \epsilon_k. \quad (15)$$

and so the right-hand side of the inequality (15) is finite. From this, we immediately get that

$$\sum_{k=1}^{\infty} \gamma_k \mathbb{E} \|\nabla \mathcal{L}(\lambda_k)\|_2^2 < \infty. \quad (16)$$

Furthermore, from (16), we get that

$$\liminf_{k \rightarrow \infty} \mathbb{E} \|\nabla \mathcal{L}(\lambda_k)\|_2^2 = 0.$$

□

## 4 Discussion

## 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>.

## Appendix

*Proof of proposition 3.9.*

$$\begin{aligned} \mathbb{E} v_Q &= \left\| \eta \sum_{i=0}^Q (I - \eta \mu \nabla_{xx}^2 F)^i \nabla_x C \right\| \leq \eta \left\| \sum_{i=0}^Q (I - \eta \mu \nabla_{xx}^2 F)^i \right\| \|\nabla_x C\| \leq \eta M \sum_{i=0}^Q (1 - \eta \mu)^i = \eta M \frac{1 - (1 - \eta \mu)^{Q+1}}{\eta \mu} \\ &= \frac{M}{\mu} (1 - (1 - \eta \mu)^{Q+1}). \end{aligned}$$

□

*Proof of Proposition 3.10.* Denote

$$\begin{aligned} A_j &= I - \eta \nabla_{xx}^2 \tilde{F}(x_k, \lambda_k, \xi_j); \\ \tilde{B} &= \nabla_x \tilde{C}(x_k, \lambda_k, \zeta); \\ A &= \mathbb{E} A_j = I - \eta \nabla_{xx}^2 F(x_k, \lambda_k); \\ B &= \mathbb{E} B = \nabla_x C(x_k, \lambda_k); \end{aligned}$$

Then

$$\begin{aligned}
\text{Var}(v_Q) &= \mathbb{E} \|v_Q - \mathbb{E}v_Q\|^2 \\
&= \mathbb{E} \left\| \eta \sum_{q=-1}^{Q-1} \prod_{j=Q-q}^Q A_j \tilde{B} - \eta \sum_{i=0}^Q A^i B \right\|^2 \\
&\leq 2\mathbb{E} \left\| \eta \sum_{q=-1}^{Q-1} \prod_{j=Q-q}^Q A_j \tilde{B} - \eta \sum_{i=0}^Q A^i \tilde{B} \right\|^2 + 2\mathbb{E} \left\| \eta \sum_{i=0}^Q A^i \tilde{B} - \eta \sum_{i=0}^Q A^i B \right\|^2 \\
&= 2\mathbb{E} \left\| \eta \sum_{q=-1}^{Q-1} \left( \prod_{j=Q-q}^Q A_j - A^{q+1} \right) \tilde{B} \right\|^2 + 2\mathbb{E} \left\| \eta \sum_{i=0}^Q A^i (\tilde{B} - B) \right\|^2.
\end{aligned}$$

Now, note that  $\mathbb{E} \left( \prod_{j=Q-q}^Q A_j - A^{q+1} \right) = 0$ , and that each  $A_i$  is independently sampled. Expanding the first term, we get

$$\begin{aligned}
&2\mathbb{E} \left\| \eta \sum_{q=-1}^{Q-1} \left( \prod_{j=Q-q}^Q A_j - A^{q+1} \right) \tilde{B} \right\|^2 + 2\mathbb{E} \left\| \eta \sum_{i=0}^Q A^i (\tilde{B} - B) \right\|^2 \\
&\leq 2\eta^2 \sum_{q=-1}^{Q-1} \mathbb{E} \left\| \left( \prod_{j=Q-q}^Q A_j - A^{q+1} \right) \tilde{B} \right\|^2 + 2\eta^2 \left\| \sum_{i=0}^Q A^i \right\|^2 \mathbb{E} \|\tilde{B} - B\|^2 \\
&\leq 2\eta^2 \sum_{q=-1}^{Q-1} \mathbb{E} \left\| \prod_{j=Q-q}^Q A_j - A^{q+1} \right\|^2 \mathbb{E} \|\tilde{B}\|^2 + 2\eta^2 \sum_{i=0}^Q \|A\|^{2i} \mathbb{E} \|\tilde{B} - B\|^2 \\
&= 2\eta^2 \sum_{q=0}^Q \mathbb{E} \left\| \prod_{j=Q-q+1}^Q A_j - A^q \right\|^2 \mathbb{E} \|\tilde{B}\|^2 + 2\eta^2 \frac{1 - \|A\|^{2Q}}{1 - \|A\|^2} \mathbb{E} \|\tilde{B} - B\|^2
\end{aligned}$$

We will now bound  $\mathbb{E}M_i$  for  $M_i = \left\| \prod_{j=Q-q+1}^Q A_j - A^q \right\|^2$ . Note, that  $M_0 = 0$ . As in the proof of proposition 3 in [2], we write

$$\prod_{j=Q-q+1}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j) = \prod_{j=Q-q+2}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j) - \eta \nabla_x^2 \tilde{F}_j \prod_{j=Q-q+2}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j)$$

Then, we have

$$\begin{aligned}
\mathbb{E}M_i &= \mathbb{E} \left\| \prod_{j=Q-i+1}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j) - [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \right\|^2 \\
&= \mathbb{E} \left\| \prod_{j=Q-i+2}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j) - \eta \nabla_x^2 \tilde{F}_j \prod_{j=Q-i+2}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j) - [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \right\|^2
\end{aligned}$$

Add and subtract  $\eta \nabla_x^2 F \prod_{j=Q-i+2}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j)$ :

$$\mathbb{E} \left\| \underbrace{\left( (I - \eta \nabla_x^2 F) \prod_{j=Q-i+2}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j) - [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \right)}_c + \underbrace{\left( (\eta \nabla_x^2 F - \eta \nabla_x^2 \tilde{F}_j) \prod_{j=Q-i+2}^Q (I - \eta \nabla_{xx}^2 \tilde{F}_j) \right)}_d \right\|^2$$

$$= \mathbb{E}\|c\|^2 + \mathbb{E}\|d\|^2 + \underbrace{2\mathbb{E}\langle c, d \rangle}_{=0 \text{ as } \mathbb{E}(\eta \nabla_x^2 F - \eta \nabla_x^2 \tilde{F}_j)=0}$$

Using convexity assumptions and Proposition 3.8, we get the bound

$$\mathbb{E}M_i \leq (1 - \eta\mu)^2 \mathbb{E}M_{i-1} + \eta^2 (1 - \eta\mu)^{2q-2} L^2.$$

Telescoping, we get

$$\mathbb{E}M_i \leq (1 - \eta\mu)^{2k} \mathbb{E}M_{i-k} + \eta^2 L^2 (1 - \eta\mu)^{2i-2} \sum_{j=1}^k (1 - \eta\mu)^{j-1}$$

Setting  $i = q, k = q$ ,

$$\mathbb{E}M_q \leq (1 - \eta\mu)^{2q-2} \mathbb{E}M_0 + \eta^2 L^2 (1 - \eta\mu)^{2q-2} \sum_{j=1}^q (1 - \eta\mu)^{j-1}$$

Note that  $\mathbb{E}(M_0) = 0$ . Thus, we finally get

$$\mathbb{E}M_q \leq \eta^2 L^2 (1 - \eta\mu)^{2q-2} \sum_{j=0}^q (1 - \eta\mu)^j = \eta^2 L^2 (1 - \eta\mu)^{2q-2} \frac{1 - (1 - \eta\mu)^{q+1}}{1 - (1 - \eta\mu)} = \frac{\eta L^2}{\mu} ((1 - \eta\mu)^{2q-2} - (1 - \eta\mu)^{3q-1}) \quad (17)$$

Hence,

$$\mathbb{E} \left\| \prod_{j=Q-q+1}^Q A_j - A^q \right\|^2 \leq \eta^2 L^2 (1 - \eta\mu)^{2q-2} \frac{1 - (1 - \eta\mu)^{q+1}}{1 - (1 - \eta\mu)} = \frac{\eta L^2}{\mu} ((1 - \eta\mu)^{2q-2} - (1 - \eta\mu)^{3q-1})$$

Hence, using the continuity and convexity assumptions, we get

$$\begin{aligned} \text{Var}(v_Q) &\leq 2\eta^2 \sum_{q=0}^Q \mathbb{E} \left\| \prod_{j=Q-q+1}^Q A_j - A^q \right\|^2 \mathbb{E} \|\tilde{B}\|^2 + 2\eta^2 \frac{1 - \|A\|^{2Q}}{1 - \|A\|^2} \mathbb{E} \|\tilde{B} - B\|^2 \\ &\leq 2\eta^2 M^2 \left( \sum_{q=0}^Q \frac{\eta L^2}{\mu} ((1 - \eta\mu)^{2q-2} - (1 - \eta\mu)^{3q-1}) + \frac{1 - (1 - \eta\mu)^{2Q}}{1 - (1 - \eta\mu)^2} \right) \\ &= \frac{2\eta^3 M^2 L^2}{\mu} \left( \frac{1 - (1 - \eta\mu)^{2Q+2}}{1 - (1 - \eta\mu)^2} - (1 - \eta\mu)^2 \frac{1 - (1 - \eta\mu)^{3Q+3}}{1 - (1 - \eta\mu)^3} \right) + 2\eta^2 M^2 \frac{1 - (1 - \eta\mu)^{2Q}}{1 - (1 - \eta\mu)^2} \\ &\leq \frac{2\eta^3 M^2 L^2}{\mu} \left( \frac{1}{1 - (1 - \eta\mu)^2} \right) + 2\eta^2 M^2 \frac{1}{1 - (1 - \eta\mu)^2} \leq \frac{2\eta^3 M^2 L^2 + 2\eta M^2}{2\mu - \eta\mu^2} \end{aligned}$$

Furthermore, since  $\eta\mu < 1$ , we have that

$$\frac{2\eta^3 M^2 L^2 + 2\eta M^2}{2\mu - \eta\mu^2} \leq \frac{2\eta^3 M^2 L^2 + 2\eta M^2}{2\mu - \mu} = 2\eta^3 M^2 L\kappa + \frac{2\eta M^2}{\mu}$$

□

*Proof of Proposition 3.11.*

$$\begin{aligned} \text{Var}(\widehat{\nabla} \mathcal{L}(\lambda_k)) &= \mathbb{E} \|\widehat{\nabla} \mathcal{L}(\lambda_k) - \mathbb{E} \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 = \\ &\mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q - (\nabla_\lambda C(x_k, \lambda_k) - \nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E} v_Q)\|^2 \\ &= \mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k)\|^2 + \mathbb{E} \|\nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E} v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q\|^2 - \end{aligned}$$

$$\begin{aligned}
& -2\mathbb{E} \left[ \left( \underbrace{\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k)}_{\mathbf{0} \text{ in expectation}} \right)^\top \left( \nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E} v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q \right) \right] \\
& = \mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k)\|^2 + \mathbb{E} \|\nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E} v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q\|^2 \\
& = \text{Var} \left( \nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) \right) + \text{Var} \left( \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q \right) \tag{18}
\end{aligned}$$

$$\begin{aligned}
& = \text{Var} \left( \nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) \right) + \text{Var} \left( \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2) \right) \text{Var} (v_Q) + \text{Var} \left( \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2) \right) \|\mathbb{E}[v_Q]\|^2 + \\
& \quad + \text{Var} (v_Q) \left\| \mathbb{E} \left[ \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2) \right] \right\|^2, \tag{19}
\end{aligned}$$

where we get from (18) to (19) using the identity  $\text{Var}[XY] = \text{Var}[X]\text{Var}[Y] + \text{Var}[X]\mathbb{E}[Y]^2 + \text{Var}[Y]\mathbb{E}[X]^2$  for independent  $X, Y$ . Now, by Proposition 3.8, and Assumption 3.4, we have

$$\text{Var} \left( \nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) \right) \leq M^2, \quad \text{Var} \left( \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2) \right) \leq L^2, \quad \left\| \mathbb{E} \left[ \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2) \right] \right\|^2 \leq K^2.$$

By Proposition 3.9, we also obtain

$$\|\mathbb{E}[v_Q]\|^2 \leq \frac{M^2}{\mu^2} (1 - (1 - \eta\mu)^{Q+1})^2 \leq \frac{M^2}{\mu^2}.$$

Finally, from Proposition 3.10, we have

$$\text{Var}(v_Q) \leq 2\eta^3 M^2 L\kappa + \frac{2\eta M^2}{\mu}.$$

Thus, we can bound (19) by

$$\left( 2\eta^3 M^2 L\kappa + \frac{2\eta M^2}{\mu} \right) (L^2 + K^2) + \kappa^2 M^2 + M^2,$$

and we are done. □