

# Stochastic Bilevel Optimization

February 20, 2023

## 1 Problem

We set up a bilevel problem. For a bounded domain  $\mathcal{D} \subset \mathbb{R}^r$ ,

$$\min_{\lambda \in \mathcal{D}} \mathcal{L}(\lambda) \triangleq C(\hat{x}(\lambda), \lambda) \quad (1)$$

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

We will denote the sampled terms as follows. For i.i.d. samples  $\zeta, \xi$ , uniformly randomly selected from the given data  $\{\zeta_i, \xi_j, i = 1, \dots, m_1; j = 1, \dots, m_2\}$ , we have:

$$\mathbb{E}_{\zeta}[\tilde{F}(x, \lambda; \zeta)] = F(x, \lambda)$$

$$\mathbb{E}_{\xi}[\tilde{C}(x, \lambda, \xi)] = C(x, \lambda)$$

For a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , will denote the conditional expectation  $\mathbb{E}_{x_k}[f(x)] = \mathbb{E}[f(x)|x = x_k]$ .

## 2 Algorithm

---

### Algorithm 1 Stochastic HOAG

---

1: At iteration  $k = 1, 2, \dots$ , choose batches of i.i.d. uniform random samples  $\mathcal{S}_{\nabla_x C}, \mathcal{S}_{\nabla_{\lambda} C}, \mathcal{S}_{\nabla_{x\lambda} F}$ , batches of i.i.d. random samples  $\mathcal{B}_j$ , stepsize  $\nu_k$ , parameter  $\eta$ , perform the following:

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

$$\|\hat{x}(\lambda_k) - x_k\| \leq \varepsilon_k$$

2.

$$v_Q = \eta \sum_{q=-1}^{Q_k-1} \prod_{j=Q_k-q}^{Q_k} \left( I - \eta \nabla_{xx}^2 \tilde{F}(x_k, \lambda_k; \mathcal{B}_j) \right) \nabla_x \tilde{C}(x_k, \lambda_k, \mathcal{S}_{\nabla_x C}), \quad (3)$$

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

$$\widehat{\nabla} \mathcal{L}(\lambda_k) = \nabla_{\lambda} \tilde{C}(x_k, \lambda_k, \mathcal{S}_{\nabla_{\lambda} C}) - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \mathcal{S}_{\nabla_{x\lambda} F})^{\top} v_Q \quad (4)$$

4. Update hyperparameters:

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


---

## 3 Stochastic HOAG

### 3.1 Assumptions

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 inner function  $F(x, \lambda)$  is  $\mu$  strongly-convex w.r.t.  $x$ . For the stochastic setting, the same assumptions hold for  $\tilde{F}(x, \lambda; \zeta)$ .*

**Assumption 3.2** (Smoothness). *Let  $(x, \lambda) \in \mathbb{R}^n \times \mathcal{D}$ . The loss function  $C(x, \lambda)$  and  $F(x, \lambda)$  satisfy the following smoothness assumptions:*

- *The function  $C(x, \lambda)$  is  $M$ -Lipschitz, i.e., for any  $(x, \lambda), (x', \lambda') \in \mathbb{R}^n \times \mathcal{D}$ ,*

$$|C(x, \lambda) - C(x', \lambda')| \leq M \|(x, \lambda) - (x', \lambda')\|.$$

- *$\nabla C(x, \lambda)$  and  $\nabla F(x, \lambda)$  are  $L$ -Lipschitz, i.e., for any  $(x, \lambda), (x', \lambda') \in \mathbb{R}^n \times \mathcal{D}$ ,*

$$\begin{aligned} \|\nabla C(x, \lambda) - \nabla C(x', \lambda')\| &\leq L \|(x, \lambda) - (x', \lambda')\|, \\ \|\nabla F(x, \lambda) - \nabla F(x', \lambda')\| &\leq L \|(x, \lambda) - (x', \lambda')\|. \end{aligned}$$

*For the stochastic case, the same assumptions hold for  $F(x, \lambda; \xi)$  and  $G(x, \lambda, \zeta)$  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  $\mathcal{L}_{\inf}$ .*

## 3.2 Preliminaries

We can obtain the exact hypergradient of the outer problem as

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

However, finding the inverse  $[\nabla_{xx} F(\hat{x}(\lambda_k), \lambda_k)]^{-1}$  can be costly, and so we consider an approximation for this term. To be more precise, consider the following:

**Lemma 3.6** (Neumann Series). *For non-singular  $A \in \mathbb{R}^{n \times n}$ ,*

$$A^{-1} = \sum_{i=0}^{\infty} (I - A)^i, \quad A \succ 0, \|A\| < 1. \quad (5)$$

Define, for i.i.d. samples  $\xi_0, \zeta_j$ , for  $j = 1, \dots, Q$ ,

$$v_Q = \eta \sum_{q=-1}^{Q-1} \prod_{j=Q-q}^Q \left( I - \eta \nabla_{xx}^2 \tilde{F}(x_k, \lambda_k; \zeta_j) \right) \nabla_x \tilde{C}(x_k, \lambda_k, \xi_0), \quad (6)$$

where we assume  $\prod_{j=Q+1}^Q (\cdot) = I$ . Choose  $\eta$ , such that  $\|I - \eta \nabla_{xx}^2 \tilde{F}(x_k, \lambda_k; \zeta_j)\| \leq (1 - \eta\mu) < 1$ . From this, with (5) we get

$$\mathbb{E}[v_Q] = \eta \sum_{i=0}^Q [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \nabla_x C(x_k, \lambda_k) \approx [\nabla_{xx}^2 F(x_k, \lambda_k)]^{-1} \nabla_x C(x_k, \lambda_k).$$

As such, denote  $\mathbb{E}[v_\infty] = [\nabla_{xx}^2 F(x_k, \lambda_k)]^{-1} \nabla_x C(x_k, \lambda_k)$ . With this in mind, and the fact that it is not generally feasible to solve the inner problem to full precision, we obtain step 3 of Algorithm 2.

### 3.3 Case I: Variable $Q_k$ , Constant batch sizes

In order to prove the required result, we could set a variable Neumann series length  $Q_k$ . Options considered:

$$\text{I)} \quad Q_k = r \log k, \quad \text{II)} \quad Q_k = r \log \frac{1}{\epsilon_k}$$

for

$$r > \frac{1}{-\log(1 - \eta\mu)}.$$

These are selected, so that the  $T_5$  term is summable. In particular, we have

$$\text{I)} \quad T_5 < \frac{\sqrt{2}LM}{\mu} k^{-1}, \quad \text{II)} \quad T_5 < \frac{\sqrt{2}LM}{\mu} \epsilon_k.$$

Note that, while the choice of  $\eta$  doesn't affect the summability of  $T_5$ , it can be problematic for the Neumann series length  $Q$ . That is, if  $\eta$  is close to 0, then  $r$  may explode. As such, in practice  $\eta$  should be chosen such that  $0 < 1 - \eta\mu < 0.7$ .

**Theorem 3.7** (Global Convergence (SGD step)). *Suppose Assumptions 3.1, 3.2, 3.3, 3.4 and 3.5 hold. In Algorithm 2, assume that the stepsize  $\nu_k$  is chosen such that*

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

*Suppose  $\eta$  is chosen, such that  $\eta\mu < 1$ . Choose  $Q_k = r \log k$ , where  $r > \frac{1}{-\log(1 - \eta\mu)}$ . Assume also, that  $\lambda_k \in \mathcal{D}$  for all  $k > 0$ . Let  $\epsilon_k$  be such that  $\|\hat{x}(\lambda_k) - x_k\| \leq \epsilon_k$ . 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*

$$\min_{k \leq K} \mathbb{E} \|\nabla \mathcal{L}(\lambda_k)\|^2 \xrightarrow{K \rightarrow \infty} 0$$

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

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

where  $\hat{\nabla} \mathcal{L}(\lambda_k)$  is the approximate hypergradient of  $\mathcal{L}$ , evaluated at  $\lambda_k$ , as defined in step 3 of Algorithm 2. Taking expectation, conditioning on  $\lambda_k$ ,

$$\begin{aligned} \mathbb{E}_{\lambda_k} [\mathcal{L}(\lambda_{k+1})] &\leq \mathcal{L}(\lambda_k) - \nu_k \nabla \mathcal{L}(\lambda_k)^\top \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] + \frac{L\nu_k^2}{2} \mathbb{E}_{\lambda_k} \left[ \left\| \hat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \right] \\ &= \mathcal{L}(\lambda_k) - \nu_k \left( \nabla \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] \right)^\top \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] - \nu_k \left\| \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] \right\|^2 + \frac{L\nu_k^2}{2} \mathbb{E}_{\lambda_k} \left[ \left\| \hat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \right] \\ &\leq \mathcal{L}(\lambda_k) + \nu_k \left\| \nabla \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] \right\| \left\| \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] \right\| - \nu_k \left\| \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] \right\|^2 + \frac{L\nu_k^2}{2} \mathbb{E}_{\lambda_k} \left[ \left\| \hat{\nabla} \mathcal{L}(\lambda_k) \right\|^2 \right] \\ &= \mathcal{L}(\lambda_k) + \nu_k \left\| \nabla \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] \right\| \left\| \mathbb{E}_{\lambda_k} [\hat{\nabla} \mathcal{L}(\lambda_k)] \right\| + \frac{L\nu_k^2}{2} \text{Var}(\hat{\nabla} \mathcal{L}(\lambda_k)) - \left( \nu_k - \frac{L\nu_k^2}{2} \right) \left\| \mathbb{E}_{\lambda_k} \hat{\nabla} \mathcal{L}(\lambda_k) \right\|^2, \end{aligned} \quad (8)$$

Let  $T_{45} = T_4\epsilon_k + T_5$ , where  $T_4, T_5$  are as in (22) of Proposition .6. From Proposition .6 and by assumption on  $\epsilon_k$ , we have that

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

Note here, that by the summability of  $\epsilon_k$ , as well as by the choice of  $Q_k$ ,  $T_{45}$  itself is summable. By reverse triangle inequality, we have that

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

Hence,

$$\|\mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k)\| \leq T_{45} + \|\nabla \mathcal{L}(\lambda_k)\| \quad \text{and} \quad \|\mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k)\| \geq \|\nabla \mathcal{L}(\lambda_k)\| - T_{45}.$$

Using the fact that for  $x > 0$ ,  $x < \max\{1, x^2\}$  we also have

$$\|\mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k)\| \leq T_{45} + 1 + \|\nabla \mathcal{L}(\lambda_k)\|^2. \quad (10)$$

Furthermore,

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

Similarly,

$$\|\mathbb{E}_{\lambda_k} \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \geq \|\nabla \mathcal{L}(\lambda_k)\|^2(1 - 2T_{45}) + T_{45}^2 - 2T_{45}. \quad (12)$$

Finally, by proposition .4, we have that

$$\text{Var}(\widehat{\nabla} \mathcal{L}(\lambda_k)) \leq 2\eta^2 M^2 \left( \frac{\kappa^2}{B(2 - \eta\mu)^2} + \frac{1}{|\mathcal{S}_{\nabla_x C}|(2\eta\mu - \eta^2\mu^2)} \right) \left( \frac{L^2}{|\mathcal{S}_{\nabla_x \lambda F}|} + K^2 \right) + \frac{\kappa^2 M^2}{|\mathcal{S}_{\nabla_x \lambda F}|} + \frac{M^2}{|\mathcal{S}_{\nabla_x \lambda C}|} \leq T_2. \quad (13)$$

Substituting (9),(10), (11), (12) and (13) into (8), and rearranging, we get

$$\begin{aligned} \left( \nu_k - \frac{L\nu_k^2}{2} \right) (\|\nabla \mathcal{L}(\lambda_k)\|^2(1 - 2T_{45}) + T_{45}^2 - 2T_{45}) - \nu_k T_{45} (\|\nabla \mathcal{L}(\lambda_k)\|^2 + 1 + T_{45}) \leq \\ \leq \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\mathcal{L}(\lambda_{k+1})] + \frac{T_2 \nu_k^2}{2L}. \end{aligned} \quad (14)$$

Collecting like terms, and further rearranging, we get

$$\|\nabla \mathcal{L}(\lambda_k)\|^2 \left( \underbrace{\nu_k - \frac{L\nu_k^2}{2} - 3\nu_k T_{45} + L\nu_k^2 T_{45}}_{\triangleq \gamma_k} \right) \leq \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\mathcal{L}(\lambda_{k+1})] + \underbrace{\frac{T_2 \nu_k^2}{2L} + \frac{L\nu_k^2 T_{45}^2}{2} + 3\nu_k T_{45} - T_{45} L\nu_k^2}_{\triangleq T_{2,k}} \quad (15)$$

Note that, by Cauchy-Schwarz,  $\sum_{k=1}^{\infty} \nu_k \epsilon_k \leq \sqrt{\sum_{k=1}^{\infty} \nu_k^2} \sqrt{\sum_{k=1}^{\infty} \epsilon_k^2} < \infty$ . Similarly, we have that  $\nu_k T_{45}$  is summable. Hence,  $T_{2,k}$  is summable. Summing (15) for  $k = 1$  to  $\infty$ , we get

$$\sum_{k=1}^{\infty} \gamma_k \|\nabla \mathcal{L}(\lambda_k)\|^2 \leq \sum_{k=1}^{\infty} (\mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} [\mathcal{L}(\lambda_{k+1})]) + \sum_{k=1}^{\infty} T_{2,k}.$$

Taking total expectation and telescoping, we get

$$\sum_{k=1}^{\infty} \gamma_k \mathbb{E} \|\nabla \mathcal{L}(\lambda_k)\|^2 \leq \mathcal{L}(\lambda_1) - \mathcal{L}_{inf} + \sum_{k=1}^{\infty} T_{2,k}, \quad (16)$$

and so the right-hand side of the inequality (16) is finite. From this, we immediately get that

$$\sum_{k=1}^{\infty} \gamma_k \mathbb{E} \|\nabla \mathcal{L}(\lambda_k)\|^2 < \infty. \quad (17)$$

Furthermore, from (17), we get that

$$\min_{k \leq K} \mathbb{E} \|\nabla \mathcal{L}(\lambda_k)\|^2 \leq \frac{\sum_{k=1}^K \gamma_k \mathbb{E} \|\nabla \mathcal{L}(\lambda_k)\|^2}{\sum_{k=1}^K \gamma_k} \xrightarrow{K \rightarrow \infty} 0$$

□

### 3.4 Case II: Variable $Q_k \mathcal{O}(1/\epsilon_k)$ batches

## 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

Here will provide the technical lemmas (with proofs) that are required for the proof of Theorem 3.7. We present an immediate consequence of the assumptions in section 3.1.

**Proposition .1** (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  $x, \lambda, \zeta$ ,*

$$\begin{aligned}\mathbb{E}_\zeta \|\nabla \tilde{C}(x, \lambda, \zeta) - \nabla C(x, \lambda)\|^2 &\leq M^2 \\ \mathbb{E}_\zeta \|\nabla_{x\lambda}^2 \tilde{F}(x, \lambda, \zeta) - \nabla_{x\lambda}^2 F(x, \lambda)\|^2 &\leq L^2 \\ \mathbb{E}_\zeta \|\nabla_{xx}^2 \tilde{F}(x, \lambda, \zeta) - \nabla_{xx}^2 F(x, \lambda)\|^2 &\leq L^2\end{aligned}$$

**Proposition .2** (Bound on  $\|\mathbb{E}v_Q\|$ ). *Suppose Assumptions 3.1, 3.2 hold. Then*

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

*Proof of proposition .2.*

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

□

We now introduce our first major result. The variance of the  $Q$ -long Neumann expansion can be controlled by the sampling of  $\nabla_{xx}^2 F$  and  $\nabla_x C$  as follows:

**Proposition .3** (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  $\eta\mu < 1$ . Let  $|\mathcal{B}_{Q-1+j}| = B(1 - \eta\mu)^{j-1}$ . Then we have that*

$$\text{Var}(v_Q) = \mathbb{E}\|v_Q - \mathbb{E}(v_Q)\|^2 \leq \frac{2\eta^2 M^2 \kappa^2}{B(1 - \eta\mu)} + \frac{2\eta M^2}{|\mathcal{S}_{\nabla_x C}|(2 - \eta\mu)\mu}.$$

*Proof of Proposition .3.* Denote

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

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

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

Then, we have (denoting  $\tilde{F}(x_k, \lambda_k; \mathcal{B}_j) = \tilde{F}_j$ )

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

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

$$\mathbb{E} \left\| \underbrace{\left( I - \eta \nabla_{xx}^2 F \right) \prod_{j=Q-i+2}^Q \left( I - \eta \nabla_{xx}^2 \tilde{F}_j \right) - [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i}_c + \underbrace{\left( \eta \nabla_{xx}^2 F - \eta \nabla_{xx}^2 \tilde{F}_{Q-i+1} \right) \prod_{j=Q-i+2}^Q \left( I - \eta \nabla_{xx}^2 \tilde{F}_j \right)}_d \right\|^2$$

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

Note, that for sample set  $\mathcal{B}_{Q-i+1} = \{\xi_j, j = 1, \dots, |\mathcal{B}_{Q-i+1}|\}$ , with i.i.d.  $\xi_j$ , by Proposition .1, we obtain

$$\begin{aligned} \mathbb{E}\|\nabla_x^2 F(x_k, \lambda_k) - \nabla_x^2 \tilde{F}_{Q-i+1}\|^2 &= \mathbb{E}\left\|\frac{1}{|\mathcal{B}_{Q-i+1}|} \sum_{j=1}^{|\mathcal{B}_{Q-i+1}|} \left(\nabla_x^2 F(x_k, \lambda_k) - \nabla_x^2 \tilde{F}(x_k, \lambda_k, \xi_j)\right)\right\|^2 \\ &= \frac{1}{|\mathcal{B}_{Q-i+1}|^2} \sum_{j=1}^{|\mathcal{B}_{Q-i+1}|} \mathbb{E}_{\xi_j} \|\nabla_x^2 F(x_k, \lambda_k) - \nabla_x^2 \tilde{F}(x_k, \lambda_k, \xi_j)\|^2 \leq \frac{L^2}{|\mathcal{B}_{Q-i+1}|} \end{aligned} \quad (18)$$

Hence, by strong convexity and Cauchy-Schwarz we obtain that

$$\mathbb{E}\|c\|^2 \leq (1 - \eta\mu)^2 \mathbb{E}M_{i-1}, \quad \mathbb{E}\|d\|^2 \leq \eta^2 (1 - \eta\mu)^{2i-2} \frac{L^2}{|\mathcal{B}_{Q-i+1}|}.$$

That is,

$$\mathbb{E}M_i \leq (1 - \eta\mu)^2 \mathbb{E}M_{i-1} + \eta^2 (1 - \eta\mu)^{2i-2} \frac{L^2}{|\mathcal{B}_{Q-i+1}|}.$$

Telescoping, we get

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

Setting  $i = q, k = q$ ,

$$\mathbb{E}M_q \leq (1 - \eta\mu)^{2q} \mathbb{E}M_0 + \eta^2 L^2 (1 - \eta\mu)^{2q-2} \sum_{j=1}^q \frac{1}{|\mathcal{B}_{Q-1+j}|}$$

Note that  $\mathbb{E}(M_0) = 0$ . Now, setting  $|\mathcal{B}_{Q-1+j}| = B(1 - \eta\mu)^{j-1}$ , we finally get

$$\begin{aligned} \mathbb{E}M_q &\leq \eta^2 L^2 (1 - \eta\mu)^{2q-2} \sum_{j=1}^q \frac{1}{B(1 - \eta\mu)^{j-1}} = \frac{\eta^2 L^2 (1 - \eta\mu)^{2q-2}}{B} \frac{\left(\frac{1}{1 - \eta\mu}\right)^q - 1}{\frac{1}{1 - \eta\mu} - 1} \\ &\leq \frac{\eta L^2 (1 - \eta\mu)^q}{B(1 - \eta\mu)\mu}. \end{aligned} \quad (19)$$

In a fashion similar to (18), we obtain  $\mathbb{E}\|\tilde{B} - B\|^2 \leq \frac{M^2}{|S_{\nabla_x C}|}$ . Hence, using the continuity and convexity assumptions, we get

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

Using the fact, that  $\sum_{q=0}^N x^q \leq \frac{1}{1-x}$ , we get

$$\text{Var}(v_Q) \leq 2\eta^2 M^2 \left( \frac{L^2}{B(1 - \eta\mu)\mu^2} + \frac{1}{|S_{\nabla_x C}|} \frac{1 - (1 - \eta\mu)^{2Q}}{1 - (1 - \eta\mu)^2} \right)$$

Finally, taking into account, that  $\kappa = \frac{L}{\mu}$ , we get

$$\text{Var}(v_Q) \leq 2\eta^2 M^2 \left( \frac{\kappa^2}{B(1 - \eta\mu)} + \frac{1}{|S_{\nabla_x C}|} \frac{1}{2\eta\mu - \eta^2\mu^2} \right)$$

□

Recall, that

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

**Proposition .4** (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  $\kappa = \frac{L}{\mu}$ . Let  $|\mathcal{B}_{Q-1+j}| = B(1 - \eta\mu)^{j-1}$ . Then the variance of the approximate hypergradient satisfies*

$$\text{Var}(\widehat{\nabla} \mathcal{L}(\lambda_k)) = \mathbb{E} \|\widehat{\nabla} \mathcal{L}(\lambda_k) - \mathbb{E} \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 \leq 2\eta^2 M^2 \left( \frac{\kappa^2}{B(2 - \eta\mu)^2} + \frac{1}{|\mathcal{S}_{\nabla_{x\lambda} C}|(2\eta\mu - \eta^2\mu^2)} \right) \left( \frac{L^2}{|\mathcal{S}_{\nabla_{x\lambda} F}|} + K^2 \right) + \frac{\kappa^2 M^2}{|\mathcal{S}_{\nabla_{x\lambda} F}|} + \frac{M^2}{|\mathcal{S}_{\nabla_\lambda C}|}.$$

*Proof of Proposition .4.*

$$\begin{aligned} \text{Var}(\widehat{\nabla} \mathcal{L}(\lambda_k)) &= \mathbb{E} \|\widehat{\nabla} \mathcal{L}(\lambda_k) - \mathbb{E} \widehat{\nabla} \mathcal{L}(\lambda_k)\|^2 = \\ &= \mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q - (\nabla_\lambda C(x_k, \lambda_k) - \nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E} v_Q)\|^2 \\ &= \mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k)\|^2 + \mathbb{E} \|\nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E} v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q\|^2 - \\ &\quad - 2\mathbb{E} \left[ \left( \underbrace{\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k)}_{\mathbf{0} \text{ in expectation}} \right)^\top \left( \nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E} v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q \right) \right] \\ &= \mathbb{E} \|\nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) - \nabla_\lambda C(x_k, \lambda_k)\|^2 + \mathbb{E} \|\nabla_{x\lambda}^2 F(x_k, \lambda_k)^\top \mathbb{E} v_Q - \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q\|^2 \\ &= \text{Var} \left( \nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) \right) + \text{Var} \left( \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2)^\top v_Q \right) \end{aligned} \quad (20)$$

$$\begin{aligned} &= \text{Var} \left( \nabla_\lambda \tilde{C}(x_k, \lambda_k, \xi_1) \right) + \text{Var} \left( \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2) \right) \text{Var} (v_Q) + \text{Var} \left( \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2) \right) \|\mathbb{E}[v_Q]\|^2 + \\ &\quad + \text{Var} (v_Q) \left\| \mathbb{E} \left[ \nabla_{x\lambda}^2 \tilde{F}(x_k, \lambda_k, \xi_2) \right] \right\|^2, \end{aligned} \quad (21)$$

where we get from (20) to (21) 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 .1, and Assumption 3.4, we have

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

By Proposition .2, 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, we have

$$\text{Var}(v_Q) \leq 2\eta^2 M^2 \left( \frac{\kappa^2}{B(2 - \eta\mu)^2} + \frac{1}{|\mathcal{S}_{\nabla_{x\lambda} C}|(2\eta\mu - \eta^2\mu^2)} \right).$$

Thus, we can bound (21) by

$$2\eta^2 M^2 \left( \frac{\kappa^2}{B(2 - \eta\mu)^2} + \frac{1}{|\mathcal{S}_{\nabla_{x\lambda} C}|(2\eta\mu - \eta^2\mu^2)} \right) \left( \frac{L^2}{|\mathcal{S}_{\nabla_{x\lambda} F}|} + K^2 \right) + \frac{\kappa^2 M^2}{|\mathcal{S}_{\nabla_{x\lambda} F}|} + \frac{M^2}{|\mathcal{S}_{\nabla_\lambda C}|},$$

and we are done.  $\square$



**Proposition .5** (Bound on  $\|\mathbb{E}v_Q - \mathbb{E}v_\infty\|$ ). *Suppose Assumptions 3.1, 3.2 hold.*

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

*Proof.* Proof of proposition .5.

$$\begin{aligned} \|\mathbb{E}v_Q - \mathbb{E}v_\infty\| &= \left\| \eta \sum_{i=0}^Q [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \nabla_x C(x_k, \lambda_k) - \eta \sum_{i=0}^{\infty} [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \nabla_x C(x_k, \lambda_k) \right\| \\ &= \left\| \eta \sum_{i=Q+1}^{\infty} [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \nabla_x C(x_k, \lambda_k) \right\| \leq \left\| \eta \sum_{i=Q+1}^{\infty} [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i \right\| \|\nabla_x C(x_k, \lambda_k)\| \\ &\leq \left\| [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^{Q+1} \right\| \left\| \underbrace{\eta \sum_{i=0}^{\infty} [I - \eta \nabla_{xx}^2 F(x_k, \lambda_k)]^i}_{[\nabla^2 F(x_k, \lambda_k)]^{-1}} \right\| M \end{aligned}$$

Note, that strong convexity of  $F$  gives us  $[\nabla^2 F]^{-1} \preceq \frac{1}{\mu} I$ , and so, we finally get

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

□

Define

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

We can then bound the difference between our approximate hypergradient  $\nabla \mathcal{L}$  and  $\mathbb{E}_{\lambda_k} \hat{\nabla} \mathcal{L}$  as follows:

**Proposition .6** (Bound on  $\|\mathbb{E}_{\lambda_k} \hat{\nabla} \mathcal{L}(\lambda_k) - \nabla \mathcal{L}(\lambda_k)\|$ . Lemma 7 of [2]). *Let*

$$T_4 = \sqrt{2} \left( L + \frac{L^2}{\mu} + \frac{M\tau}{\mu} + \frac{LM\rho}{\mu^2} \right), \quad T_5 = \sqrt{2} \frac{LM(1 - \eta\mu)^Q}{\mu}. \quad (22)$$

*Then we have that*

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

*where  $\epsilon_k$  is such that  $\|\hat{x}(\lambda_k) - x_k\| \leq \epsilon_k$ .*

The above result is based on Proposition .5 and the fact that

$$\|\mathbb{E}_{\lambda_k} \hat{\nabla} \mathcal{L}(\lambda_k) - \nabla \mathcal{L}(\lambda_k)\|^2 \leq 2\|\tilde{\nabla} \mathcal{L}(\lambda_k) - \nabla \mathcal{L}(\lambda_k)\|^2 + 2\|\tilde{\nabla} \mathcal{L}(\lambda_k) - \mathbb{E}_{\lambda_k} \hat{\nabla} \mathcal{L}(\lambda_k)\|^2.$$