

# A UNIFIED THEORY OF ERROR FEEDBACK AND VARIANCE REDUCTION FOR NON-CONVEX OPTIMIZATION

**Kai Yi**

King Abdullah University of Science and Technology (KAUST), Thuwal, Saudi Arabia  
kai.yi@kaust.edu.sa

## ABSTRACT

In this project, we extend EF-BV to non-convex setting.

## CONTENTS

### A Appendix

2

## REFERENCES

Laurent Condat, Kai Yi, and Peter Richtárik. Ef-bv: A unified theory of error feedback and variance reduction mechanisms for biased and unbiased compression in distributed optimization. [arXiv preprint arXiv:2205.04180](#), 2022.

Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. EF21: A new, simpler, theoretically better, and practically faster error feedback. In [Proc. of 35th Conf. Neural Information Processing Systems \(NeurIPS\)](#), 2021.

## A APPENDIX

$$g_i^t = h_i^t + \nu \mathcal{C}(\nabla f_i(x^t) - h_i^t); \quad h_i^{t+1} = h_i^t + \lambda \mathcal{C}(\nabla f_i(x^t) - h_i^t)$$

*Proof.* We first bound  $\mathbb{E}[\|g_i^{t+1} - \nabla f_i(x^{t+1})\|^2]$ .

Given any compressor  $\mathcal{C}$ , we can do a *bias-variance decomposition* of the compression error. That means, for every  $x \in \mathbb{R}^d$ ,

$$\mathbb{E}[\|\mathcal{C}(x) - x\|^2] = \underbrace{\|\mathbb{E}[\mathcal{C}(x)] - x\|^2}_{\text{bias}} + \underbrace{\mathbb{E}[\|\mathcal{C}(x) - \mathbb{E}[\mathcal{C}(x)]\|^2]}_{\text{variance}}. \quad (1)$$

For every  $t \geq 0$ , conditionally on  $x^t, h^t$  and  $(h_i^t)_{i=1}^n$ ,

$$\begin{aligned} \mathbb{E}[\|g^{t+1} - \nabla f(x^t)\|^2] &= \mathbb{E}\left[\left\|\frac{1}{n} \sum_{i=1}^n \left(h_i^t - \nabla f_i(x^t) + \nu \mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)\right)\right\|^2\right] \\ &\stackrel{(1)}{=} \left\|\frac{1}{n} \sum_{i=1}^n \left(h_i^t - \nabla f_i(x^t) + \nu \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\right)\right\|^2 \\ &\quad + \nu^2 \mathbb{E}\left[\left\|\frac{1}{n} \sum_{i=1}^n \left(\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t) - \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\right)\right\|^2\right] \\ &\leq \left\|\frac{1}{n} \sum_{i=1}^n \left(h_i^t - \nabla f_i(x^t) + \nu \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\right)\right\|^2 \\ &\quad + \nu^2 \frac{\omega_{\text{av}}}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - h_i^t\|^2, \end{aligned}$$

where the last inequality follows from Eqn. 7 in Condat et al. (2022). In addition,

$$\begin{aligned} &\left\|\frac{1}{n} \sum_{i=1}^n \left(h_i^t - \nabla f_i(x^t) + \nu \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\right)\right\| \\ &\leq \left\|\frac{1}{n} \sum_{i=1}^n \left(\nu(h_i^t - \nabla f_i(x^t)) + \nu \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\right)\right\| \\ &\quad + (1 - \nu) \left\|\frac{1}{n} \sum_{i=1}^n (h_i^t - \nabla f_i(x^t))\right\| \\ &\leq \frac{\nu}{n} \sum_{i=1}^n \|h_i^t - \nabla f_i(x^t) + \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\| \\ &\quad + \frac{1 - \nu}{n} \sum_{i=1}^n \|h_i^t - \nabla f_i(x^t)\| \\ &\leq \frac{\nu\eta}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - h_i^t\| + \frac{1 - \nu}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - h_i^t\| \\ &= \frac{1 - \nu + \nu\eta}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - h_i^t\|. \end{aligned}$$

where the last step is obtained by using the definition of the general compressor.

Therefore,

$$\left\| \frac{1}{n} \sum_{i=1}^n \left( h_i^t - \nabla f_i(x^t) + \nu \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)] \right) \right\|^2 \leq \frac{(1 - \nu + \nu\eta)^2}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - h_i^t\|^2,$$

and, conditionally on  $x^t, h^t$  and  $(h_i^t)_{i=1}^n$ ,

$$\mathbb{E}[\|g^{t+1} - \nabla f(x^t)\|^2] \leq ((1 - \nu + \nu\eta)^2 + \nu^2 \omega_{\text{av}}) \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^t) - h_i^t\|^2.$$

Using the Tower property and we obtain the unconditioned term,

$$\mathbb{E}[\|g^{t+1} - \nabla f(x^t)\|^2] \leq ((1 - \nu + \nu\eta)^2 + \nu^2 \omega_{\text{av}}) \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|\nabla f_i(x^t) - h_i^t\|^2].$$

Since we have the property (Richtárik et al., 2021, Lemma 4), for every  $t \geq 0$ ,

$$f(x^{t+1}) - f^{\text{inf}} \leq f(x^t) - f^{\text{inf}} - \frac{\gamma}{2} \|\nabla f(x^t)\|^2 + \frac{\gamma}{2} \|g^{t+1} - \nabla f(x^t)\|^2 + \left( \frac{L}{2} - \frac{1}{2\gamma} \right) \|x^{t+1} - x^t\|^2.$$

Thus, for every  $t \geq 0$ , conditionally on  $x^t, h^t$  and  $(h_i^t)_{i=1}^n$ ,

$$\begin{aligned} \mathbb{E}[f(x^{t+1}) - f^{\text{inf}}] &\leq \mathbb{E}[f(x^t) - f^{\text{inf}}] - \frac{\gamma}{2} \mathbb{E}[\|\nabla f(x^t)\|^2] \\ &\quad + \left( \frac{L}{2} - \frac{1}{2\gamma} \right) \mathbb{E}[\|x^{t+1} - x^t\|^2] + \frac{\gamma}{2} \mathbb{E}[\|g^{t+1} - \nabla f(x^t)\|^2]. \end{aligned}$$

Now, let us study the control variates  $h_i^t$ . Let  $s > 0$ . Using the Peter–Paul inequality  $\|a + b\|^2 \leq (1 + s)\|a\|^2 + (1 + s^{-1})\|b\|^2$ , for any vectors  $a$  and  $b$ , we have, for every  $t \geq 0$  and  $i \in \mathcal{I}_n$ ,

$$\begin{aligned} \|\nabla f_i(x^{t+1}) - h_i^{t+1}\|^2 &= \|h_i^t - \nabla f_i(x^{t+1}) + \lambda \mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)\|^2 \\ &\leq (1 + s) \|h_i^t - \nabla f_i(x^t) + \lambda \mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)\|^2 \\ &\quad + (1 + s^{-1}) \|\nabla f_i(x^{t+1}) - \nabla f_i(x^t)\|^2 \\ &\leq (1 + s) \|h_i^t - \nabla f_i(x^t) + \lambda \mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)\|^2 \\ &\quad + (1 + s^{-1}) L_i^2 \|x^{t+1} - x^t\|^2. \end{aligned}$$

Moreover, conditionally on  $x^t, h^t$  and  $(h_i^t)_{i=1}^n$ ,

$$\begin{aligned} \mathbb{E}[\|h_i^t - \nabla f_i(x^t) + \lambda \mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)\|^2] &= \|h_i^t - \nabla f_i(x^t) + \lambda \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\|^2 \\ &\quad + \lambda^2 \mathbb{E}[\|\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t) - \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\|^2] \\ &\leq \|h_i^t - \nabla f_i(x^t) + \lambda \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\|^2 \\ &\quad + \lambda^2 \omega \|\nabla f_i(x^t) - h_i^t\|^2. \end{aligned}$$

In addition,

$$\begin{aligned} \|h_i^t - \nabla f_i(x^t) + \lambda \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\| &\leq \|\lambda(h_i^t - \nabla f_i(x^t)) + \lambda \mathbb{E}[\mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)]\| \\ &\quad + (1 - \lambda) \|h_i^t - \nabla f_i(x^t)\| \\ &\leq \lambda \eta \|h_i^t - \nabla f_i(x^t)\| + (1 - \lambda) \|\nabla f_i(x^t) - h_i^t\| \\ &= (1 - \lambda + \lambda \eta) \|\nabla f_i(x^t) - h_i^t\|. \end{aligned}$$

Therefore, conditionally on  $x^t, h^t$  and  $(h_i^t)_{i=1}^n$ ,

$$\mathbb{E}[\|h_i^t - \nabla f_i(x^t) + \lambda \mathcal{C}_i^t(\nabla f_i(x^t) - h_i^t)\|^2] \leq ((1 - \lambda + \lambda \eta)^2 + \lambda^2 \omega) \|\nabla f_i(x^t) - h_i^t\|^2$$

and

$$\begin{aligned} \mathbb{E} \left[ \left\| \nabla f_i(x^{t+1}) - h_i^{t+1} \right\|^2 \right] &\leq (1+s)((1-\lambda+\lambda\eta)^2 + \lambda^2\omega) \left\| \nabla f_i(x^t) - h_i^t \right\|^2 \\ &\quad + (1+s^{-1})L_i^2 \mathbb{E} \left[ \left\| x^{t+1} - x^t \right\|^2 \right], \end{aligned}$$

so that

$$\begin{aligned} \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \left\| \nabla f_i(x^{t+1}) - h_i^{t+1} \right\|^2 \right] &\leq (1+s)((1-\lambda+\lambda\eta)^2 + \lambda^2\omega) \frac{1}{n} \sum_{i=1}^n \left\| \nabla f_i(x^t) - h_i^t \right\|^2 \\ &\quad + (1+s^{-1})\tilde{L}^2 \mathbb{E} \left[ \left\| x^{t+1} - x^t \right\|^2 \right]. \end{aligned}$$

Let  $\theta > 0$ ; its value will be set to  $\theta^*$  later on. We introduce the Lyapunov function, for every  $t \geq 0$ ,

$$\Psi^t := f(x^t) - f^{\inf} + \frac{\gamma}{2\theta} \frac{1}{n} \sum_{i=1}^n \left\| \nabla f_i(x^t) - h_i^t \right\|^2.$$

Hence, for every  $t \geq 0$ , conditionally on  $x^t, h^t$  and  $(h_i^t)_{i=1}^n$ , we have

$$\begin{aligned} \mathbb{E}[\Psi^{t+1}] &\leq \mathbb{E}[f(x^t) - f^{\inf}] - \frac{\gamma}{2} \mathbb{E} \left[ \left\| \nabla f(x^t) \right\|^2 \right] \\ &\quad + \frac{\gamma}{2\theta} \left( \theta((1-\nu+\nu\eta)^2 + \nu^2\omega_{\text{av}}) + (1+s)((1-\lambda+\lambda\eta)^2 + \lambda^2\omega) \right) \frac{1}{n} \sum_{i=1}^n \left\| \nabla f_i(x^t) - h_i^t \right\|^2 \\ &\quad + \left( \frac{L}{2} - \frac{1}{2\gamma} + \frac{\gamma}{2\theta} (1+s^{-1})\tilde{L}^2 \right) \mathbb{E} \left[ \left\| x^{t+1} - x^t \right\|^2 \right]. \end{aligned}$$

Making use of  $r$  and  $r_{\text{av}}$  and setting  $\theta = s(1+s)\frac{r}{r_{\text{av}}}$ , we can rewrite the above equation as:

$$\begin{aligned} \mathbb{E}[\Psi^{t+1}] &\leq \mathbb{E}[f(x^t) - f^{\inf}] - \frac{\gamma}{2} \mathbb{E} \left[ \left\| \nabla f(x^t) \right\|^2 \right] + \frac{\gamma}{2\theta} \left( \theta r_{\text{av}} + (1+s)r \right) \frac{1}{n} \sum_{i=1}^n \left\| \nabla f_i(x^t) - h_i^t \right\|^2 \\ &\quad + \left( \frac{L}{2} - \frac{1}{2\gamma} + \frac{\gamma}{2\theta} (1+s^{-1})\tilde{L}^2 \right) \mathbb{E} \left[ \left\| x^{t+1} - x^t \right\|^2 \right] \\ &= \mathbb{E}[f(x^t) - f^{\inf}] - \frac{\gamma}{2} \mathbb{E} \left[ \left\| \nabla f(x^t) \right\|^2 \right] + \frac{\gamma}{2\theta} (1+s)^2 \frac{r}{n} \sum_{i=1}^n \left\| \nabla f_i(x^t) - h_i^t \right\|^2 \\ &\quad + \left( \frac{L}{2} - \frac{1}{2\gamma} + \frac{\gamma}{2s^2} \frac{r_{\text{av}}}{r} \tilde{L}^2 \right) \mathbb{E} \left[ \left\| x^{t+1} - x^t \right\|^2 \right]. \end{aligned}$$

We now choose  $\gamma$  small enough so that

$$L - \frac{1}{\gamma} + \frac{\gamma}{s^2} \frac{r_{\text{av}}}{r} \tilde{L}^2 \leq 0. \quad (2)$$

A sufficient condition for equation 2 to hold is (Richtárik et al., 2021, Lemma 5):

$$0 < \gamma \leq \frac{1}{L + \tilde{L} \sqrt{\frac{r_{\text{av}}}{r} \frac{1}{s}}}. \quad (3)$$

Then, assuming that equation 3 holds, we have, for every  $t \geq 0$ , conditionally on  $x^t, h^t$  and  $(h_i^t)_{i=1}^n$ ,

$$\begin{aligned} \mathbb{E}[\Psi^{t+1}] &\leq \mathbb{E}[f(x^t) - f^{\inf}] - \frac{\gamma}{2} \mathbb{E} \left[ \left\| \nabla f(x^t) \right\|^2 \right] + \frac{\gamma}{2\theta} (1+s)^2 \frac{r}{n} \sum_{i=1}^n \left\| \nabla f_i(x^t) - h_i^t \right\|^2 \\ &\leq \max(1 - \gamma\mu, (1+s)^2 r) \Psi^t. \end{aligned}$$

We see that  $s$  must be small enough so that  $(1+s)^2 r < 1$ ; this is the case with  $s = s^*$ , so that  $(1+s^*)^2 r = \frac{r+1}{2} < 1$ . Therefore, we set  $s = s^*$ , and, accordingly,  $\theta = \theta^*$ . Then, for every  $t \geq 0$ , conditionally on  $x^t, h^t$  and  $(h_i^t)_{i=1}^n$ ,

$$\mathbb{E}[\Psi^{t+1}] \leq \max\left(1 - \gamma\mu, \frac{r+1}{2}\right) \Psi^t.$$

□