

Theorem. Let v_i 's be independent Rademacher variables and $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth. Then, we have the following bounds for the control variate forward gradient $\mathbf{h}_{\mathbf{v},\epsilon}(\boldsymbol{\theta})$ for all $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\|\mathbb{E}[\mathbf{h}_{\mathbf{v},\epsilon}(\boldsymbol{\theta})] - \nabla f(\boldsymbol{\theta})\| \leq \frac{\epsilon L d^{3/2}}{2}, \quad (1)$$

and

$$\mathbb{E}[\|\mathbf{h}_{\mathbf{v},\epsilon}(\boldsymbol{\theta}) - \nabla f(\boldsymbol{\theta})\|^2] \leq \frac{\epsilon^2 L^2 d^3}{4} + 2(d-1) \frac{\epsilon L d^{3/2}}{2} \|\nabla f(\boldsymbol{\theta}) - \nabla \hat{f}(\boldsymbol{\theta})\| + (d-1) \|\nabla f(\boldsymbol{\theta}) - \nabla \hat{f}(\boldsymbol{\theta})\|^2 \quad (2)$$

Proof. Since f is L -smooth, we have

$$|f(\mathbf{x}) - f(\mathbf{y}) - \nabla f(\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y})| \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d). \quad (3)$$

This inequality with $\boldsymbol{\theta} + \epsilon \mathbf{v}$ and $\boldsymbol{\theta}$ yields

$$|f(\boldsymbol{\theta} + \epsilon \mathbf{v}) - f(\boldsymbol{\theta}) - \nabla f(\boldsymbol{\theta}) \cdot \epsilon \mathbf{v}| \leq \frac{L}{2} \|\epsilon \mathbf{v}\|^2 \quad (4)$$

$$= \frac{\epsilon^2 L d}{2}. \quad (5)$$

We denote $D_{\mathbf{v},\epsilon}(\boldsymbol{\theta}) := \frac{f(\boldsymbol{\theta} + \epsilon \mathbf{v}) - f(\boldsymbol{\theta})}{\epsilon} - \nabla f(\boldsymbol{\theta}) \cdot \mathbf{v}$. Note that $|D_{\mathbf{v},\epsilon}(\boldsymbol{\theta})| \leq \epsilon L d/2$. Using this inequality, we can show inequality (1) as follows:

$$\|\mathbb{E}[\mathbf{h}_{\mathbf{v},\epsilon}(\boldsymbol{\theta})] - \nabla f(\boldsymbol{\theta})\| = \|\mathbb{E}[\mathbf{g}_{\mathbf{v},\epsilon}(\boldsymbol{\theta}) - \hat{\mathbf{g}}_{\mathbf{v}}(\boldsymbol{\theta}) + \mathbb{E}[\hat{\mathbf{g}}_{\mathbf{v}}(\boldsymbol{\theta})]] - \mathbb{E}[\mathbf{g}_{\mathbf{v}}(\boldsymbol{\theta})]\| \quad (6)$$

$$= \|\mathbb{E}[\mathbf{g}_{\mathbf{v},\epsilon}(\boldsymbol{\theta}) - \mathbf{g}_{\mathbf{v}}(\boldsymbol{\theta})]\| \quad (7)$$

$$\leq \mathbb{E}[\|\mathbf{g}_{\mathbf{v},\epsilon}(\boldsymbol{\theta}) - \mathbf{g}_{\mathbf{v}}(\boldsymbol{\theta})\|] \quad (8)$$

$$= \mathbb{E}\left[\left\|\frac{f(\boldsymbol{\theta} + \epsilon \mathbf{v}) - f(\boldsymbol{\theta})}{\epsilon} \mathbf{v} - (\nabla f(\boldsymbol{\theta}) \cdot \mathbf{v}) \mathbf{v}\right\|\right] \quad (9)$$

$$= \mathbb{E}[|D_{\mathbf{v},\epsilon}(\boldsymbol{\theta})| \|\mathbf{v}\|] \quad (10)$$

$$\leq \frac{\epsilon L d}{2} d^{1/2} \quad (11)$$

$$= \frac{\epsilon L d^{3/2}}{2}. \quad (12)$$

Next, we show inequality (2). To simplify the notation, we omit the argument θ . Then, we have

$$\mathbb{E}[\|\mathbf{h}_{\mathbf{v},\epsilon} - \nabla f\|^2] = \mathbb{E}[\|(\mathbf{g}_{\mathbf{v},\epsilon} - \mathbf{g}_{\mathbf{v}}) + (\mathbf{h}_{\mathbf{v}} - \nabla f)\|^2] \quad (13)$$

$$= \mathbb{E}[\|D_{\mathbf{v},\epsilon}\mathbf{v} + (\mathbf{h}_{\mathbf{v}} - \nabla f)\|^2] \quad (14)$$

$$= \mathbb{E}[\|D_{\mathbf{v},\epsilon}\mathbf{v}\|^2] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot (\mathbf{h}_{\mathbf{v}} - \nabla f)] + \mathbb{E}[\|\mathbf{h}_{\mathbf{v}} - \nabla f\|^2] \quad (15)$$

$$= \mathbb{E}[D_{\mathbf{v},\epsilon}^2 \|\mathbf{v}\|^2] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot (\mathbf{g}_{\mathbf{v}} - \hat{\mathbf{g}}_{\mathbf{v}} + \nabla \hat{f} - \nabla f)] + (d-1) \|\nabla f - \nabla \hat{f}\|^2 \quad (16)$$

$$= \mathbb{E}[D_{\mathbf{v},\epsilon}^2 \|\mathbf{v}\|^2] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot (\mathbf{g}_{\mathbf{v}} - \hat{\mathbf{g}}_{\mathbf{v}})] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot (\nabla \hat{f} - \nabla f)] + (d-1) \|\nabla f - \nabla \hat{f}\|^2 \quad (17)$$

$$= \mathbb{E}[D_{\mathbf{v},\epsilon}^2 \|\mathbf{v}\|^2] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot ((\nabla f - \nabla \hat{f}) \cdot \mathbf{v})\mathbf{v}] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot (\nabla \hat{f} - \nabla f)] + (d-1) \|\nabla f - \nabla \hat{f}\|^2 \quad (18)$$

$$= \mathbb{E}[D_{\mathbf{v},\epsilon}^2 \|\mathbf{v}\|^2] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}((\nabla f - \nabla \hat{f}) \cdot \mathbf{v}) \|\mathbf{v}\|^2] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot (\nabla \hat{f} - \nabla f)] + (d-1) \|\nabla f - \nabla \hat{f}\|^2 \quad (19)$$

$$= d\mathbb{E}[D_{\mathbf{v},\epsilon}^2] + 2d\mathbb{E}[D_{\mathbf{v},\epsilon}((\nabla f - \nabla \hat{f}) \cdot \mathbf{v})] + 2\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot (\nabla \hat{f} - \nabla f)] + (d-1) \|\nabla f - \nabla \hat{f}\|^2 \quad (20)$$

$$= d\mathbb{E}[D_{\mathbf{v},\epsilon}^2] + 2(d-1)\mathbb{E}[D_{\mathbf{v},\epsilon}\mathbf{v} \cdot (\nabla f - \nabla \hat{f})] + (d-1) \|\nabla f - \nabla \hat{f}\|^2 \quad (21)$$

$$\leq d\left(\frac{\epsilon L d}{2}\right)^2 + 2(d-1)\frac{\epsilon L d^{3/2}}{2} \|\nabla f - \nabla \hat{f}\| + (d-1) \|\nabla f - \nabla \hat{f}\|^2 \quad (22)$$

$$= \frac{\epsilon^2 L^2 d^3}{4} + 2(d-1)\frac{\epsilon L d^{3/2}}{2} \|\nabla f - \nabla \hat{f}\| + (d-1) \|\nabla f - \nabla \hat{f}\|^2. \quad (23)$$

□