

Variants of zeroth-order
gradient estimator with
improved bounds

Previously, we assumed

$$f \in \mathcal{C}^3 \text{ (3-times cont diffble)}$$

& showed

$$\mathbb{E} \|\hat{\nabla} f(x) - \nabla f(x)\| \leq C_d^2$$

Let's fix the unified estimate

Recall

$$\begin{aligned} \nabla f(x) &= \left[\frac{f(x+\delta u) - f(x-\delta u)}{2\delta} \right] v \\ &+ \left[\frac{\xi^+ - \xi^-}{2\delta} \right] v \end{aligned}$$

Assume $\mathbb{E}[\xi^+ - \xi^- | v] = 0$
 $\mathbb{E}[v v^T] = I$

Variant 1: f is convex + L -smooth

$$\frac{\delta \nabla f(x)^T u}{2\delta} \stackrel{\text{convex}}{\leq} \frac{f(x+\delta u) - f(x)}{2\delta} \stackrel{L\text{-smooth}}{\leq} \frac{\delta \nabla f(x)^T u + \frac{L}{2} \delta^2 \|u\|^2}{2\delta}$$

A similar inequality holds for $f(x-\delta u)$

$$\frac{\nabla f(x)^T u - \frac{L}{2} \delta \|u\|^2}{2} \leq \frac{f(x+\delta u) - f(x-\delta u)}{2\delta} \leq \frac{\nabla f(x)^T u + \frac{L}{2} \delta \|u\|^2}{2}$$

$$\text{Let } \phi(x, \delta, u) = \frac{1}{\delta} \left[\frac{f(x+\delta u) - f(x-\delta u)}{2\delta} - \nabla f(x)^T u \right]$$

$$\text{Then, } |\phi(x, \delta, u)| \leq \frac{L \|u\|^2}{2}$$

$$E \hat{\nabla} f(x) = E \left[\nabla \left(\frac{f(x+\delta u) - f(x-\delta u)}{2\delta} \right) \right]$$

$$= E \left[\nabla u^T \nabla f(x) + \delta \phi(x, \delta, u) \nabla \right]$$

Using $E(\nabla u^T) = I$,

$$E \|\hat{\nabla} f(x) - \nabla f(x)\| \leq \delta E \|\phi(\cdot) \nabla\|$$

$$\leq \frac{\delta L E \|\nabla u^2\|}{2}$$

Assume $E \|\nabla u^2\| \leq C_4$

Then,

$$E \|\hat{\nabla} f(x) - \nabla f(x)\| \leq \text{const} \times \delta$$

For the claim,

$$E \|\hat{\nabla} f(x) - E[\hat{\nabla} f(x) | x]\|^2 \leq \frac{\text{const}}{\delta^2},$$

The proof is as in the case $f \in \mathbb{C}^3$.

Gaussian Smoothing

"Nesterov Spokoiny, 2017"

Goes back to "Kulkarni-Kulchitsky, 1972",

See also Kushner-Clark, 1978

Aka "Gaussian smoothed functional".

$$\hat{\nabla} f(x) = \Delta \left(\frac{(f(x + \delta \Delta) + \xi^+) - (f(x) + \xi)}{\delta} \right)$$

$\Delta \rightarrow$ vector of $N(0, 1)$

$\Delta \sim N(0, I_d)$

$$\mathbb{E}(\xi^+) = \mathbb{E}(\xi) = 0$$

Then,

$$\|E \hat{\nabla} f(x) - \nabla f(x)\| \leq C_1 \delta$$

Pf^o

||

"Smoothed functional"

$$f_{\delta}(x) = \frac{1}{(2\pi)^{d/2}} \int_{-\infty}^{\infty} f(x + \delta u) \exp\left(-\frac{\|u\|^2}{2}\right) du$$

Next lecture

Lecture - 19

Gaussian smoothing: $x \in \mathbb{R}^d$

$$f_\delta(x) = \frac{1}{(2\pi)^{d/2}} \int f(x + \delta u) e^{-\frac{1}{2} \|u\|^2} du$$

why?

$$= \frac{1}{(2\pi)^{d/2}} \int_{\mathbb{R}^d} f(y) e^{-\frac{1}{2\delta^2} \|y-x\|^2} dy$$

$y = x + \delta u$ Jacobian matrix = $\begin{bmatrix} \delta & & 0 \\ & \ddots & \\ 0 & & \delta \end{bmatrix}$

$|Jacobian| = \delta^d$

$\frac{dy}{\delta^d} = du$

$$\nabla f_\delta(x) = \frac{1}{(2\pi)^{d/2}} \delta^{d+2} \int f(y) e^{-\frac{1}{2\delta^2} \|y-x\|^2} (y-x) dy$$

$$= \frac{1}{(2\pi)^{d/2}} \delta \int f(x + \delta u) e^{-\frac{1}{2} \|u\|^2} u du$$

$$\nabla f_\delta(x) = \frac{1}{(2\pi)^{d/2}} \int \frac{f(x + \delta u) - f(x)}{\delta} e^{-\frac{1}{2} \|u\|^2} u du$$

used

$$\int u e^{-\frac{1}{2}\|u\|^2} du = 0$$

Also, $\nabla f_\delta(x) = \frac{1}{(2\pi)^{d/2}} \int \frac{f(x+du) - f(x-du)}{2\delta} e^{-\frac{1}{2}\|u\|^2} u du$

Now,

$$\|\nabla f_\delta(x) - \nabla f(x)\|$$

$$\leq \frac{1}{(2\pi)^{d/2}} \delta \int |f(x+du) - f(x) - \delta \langle \nabla f(x), u \rangle| \|u\| e^{-\frac{1}{2}\|u\|^2} du$$

used

$$\frac{1}{(2\pi)^{d/2}} \int \langle \nabla f(x), u \rangle e^{-\frac{1}{2}\|u\|^2} u du = \nabla f(x)$$

Since $\frac{1}{(2\pi)^{d/2}} \int u u^T e^{-\frac{1}{2}\|u\|^2} du = I$

$$\leq \frac{1}{(2\pi)^{d/2}} \times \frac{\delta L}{2} \int \|u\|^3 e^{-\frac{1}{2}\|u\|^2} du$$

by $|f(y) - f(x) - \langle \nabla f(x), y-x \rangle| \leq \frac{1}{2} L \|x-y\|^2$

$$\leq \frac{\delta L}{2} (d+3)^{3/2}$$

↑

$$\text{wird } \frac{1}{(2\pi)^{d/2}} \int \|u\|^3 e^{-\frac{1}{2}\|u\|^2} du = (d+3)^{3/2}$$

$$\nabla f_\delta(x) = \frac{1}{(2\pi)^{d/2}} \delta^{d+2} \int f(y) e^{-\frac{1}{2\delta^2}\|y-x\|^2} (y-x) dy$$

$$= \frac{1}{(2\pi)^{d/2}} \delta \int f(x+\delta u) e^{-\frac{1}{2}\|u\|^2} u du$$

$$\nabla f_\delta(x) = \frac{1}{\delta c} \int f(x+\delta u) e^{-\frac{1}{2}\|u\|^2} u du$$

$$u = -v, \quad du = -dv$$

$$\nabla f_\delta(x) = \frac{1}{\delta c} \int f(x-\delta v) e^{-\frac{1}{2}\|v\|^2} v dv$$

$$\nabla f_{\delta}(x) = \frac{1}{(2\pi)^{d/2}} \int \frac{f(x+\delta u) - f(x)}{\delta} e^{-\frac{1}{2}\|u\|^2} u \, du$$

Now? \rightarrow

$$\frac{1}{(2\pi)^{d/2}} \int \frac{f(x) - f(x-\delta u)}{\delta} e^{-\frac{1}{2}\|u\|^2} u \, du$$

$x + \delta u = \vartheta$, $x = \vartheta - \delta u$

$$\nabla f_{\delta}(x) = \frac{1}{(2\pi)^{d/2}} \int \frac{f(\vartheta) - f(\vartheta - \delta u)}{\delta} e^{-\frac{1}{2}\|u\|^2} u \, du$$

A Calculation:

$$\frac{1}{(2\pi)^{d/2}} \int_{-\infty}^{\infty} \nabla f(x)^{\top} u \exp\left(-\frac{\|u\|^2}{2}\right) u \, du$$

$$= \sum_{i=1}^d \nabla_i f(x) \frac{1}{(2\pi)^{d/2}} \int_{-\infty}^{\infty} u_i \exp\left(-\frac{\|u\|^2}{2}\right) u \, du$$

$$= \sum_{i=1}^d \nabla_i f(x) \frac{1}{(2\pi)^{d/2}}$$

$$\times \int_{-\infty}^{\infty} (u_1 u_i, u_2 u_i, \dots, u_{i-1} u_i, u_{i+1} u_i, \dots, u_i u_d) \times \exp\left(-\frac{\|u\|^2}{2}\right) du$$

$$= \sum_{i=1}^d \nabla_i f(x) \frac{1}{(2\pi)^{d/2}} \int u_i^2 \exp\left(-\frac{\|u\|^2}{2}\right) du$$

(wed $\int_{i \neq j} u_i u_j \exp\left(-\frac{\|u\|^2}{2}\right) du = 0$)

$$= \sum_{i=1}^d \nabla_i f(x) \frac{1}{(2\pi)^{d/2}} \left(\prod_{j \neq i} \int \exp\left(-\frac{u_j^2}{2}\right) du_j \right) \times \int u_i^2 \exp\left(-\frac{u_i^2}{2}\right) du_i$$

$$= \nabla f(x)$$

Main Point?

$$\nabla f_{\delta}(x) = \frac{1}{(2\pi)^{d/2}} \int \frac{f(x+\delta u) - f(x)}{\delta} e^{-\frac{1}{2}\|u\|^2} u \, du$$

$$= \mathbb{E} \left(\left(\frac{f(x+\delta u) - f(x)}{\delta} \right) u \right)$$

An estimate of $\nabla f_{\delta}(x)$ is $u_1 \sim \mathcal{N}(0, I_2)$

$$\hat{g} = \frac{(f(x+\delta u_1) - f(x)) u_1}{\delta}$$

$$\mathbb{E} \hat{g} = \nabla f_{\delta}(x)$$

$$\|\nabla f_{\delta} - \nabla f\| \leq (\text{const}) \times \delta$$

Bias/Variance guarantees:

See Sec 3.3.4 of the book

Assume f is convex & F is L -smooth.

Then, $\| \mathbb{E} \hat{\nabla} f(x) - \nabla f(x) \| \leq C_1 \delta$ and

$$\mathbb{E} \| \hat{\nabla} f(x) - \mathbb{E} \hat{\nabla} f(x) \|^2 \leq C_2 + C_3 \delta^2$$

Proof. As in the proof of Proposition 3.2, for any convex function h with an L -Lipschitz gradient, for any $\delta > 0$, we have

$$\frac{\langle \nabla h(\theta), \delta u \rangle}{2\delta} \leq \frac{h(\theta + \delta u) - h(\theta)}{2\delta} \leq \frac{\langle \nabla h(\theta), \delta u \rangle + (L/2) \|\delta u\|^2}{2\delta}.$$

Using similar inequalities for $h(\theta - \delta u)$, we obtain

$$\langle \nabla h(\theta), u \rangle - \frac{L\delta \|u\|^2}{2} \leq \frac{h(\theta + \delta u) - h(\theta - \delta u)}{2\delta} \leq \langle \nabla h(\theta), u \rangle + \frac{L\delta \|u\|^2}{2}.$$

Letting $\phi(\theta, \delta, u) := \frac{1}{\delta} \left(\frac{h(\theta + \delta u) - h(\theta - \delta u)}{2\delta} - \langle \nabla h(\theta), u \rangle \right)$, we get

$$|\phi(\theta, \delta, u)| \leq \frac{L}{2} \|u\|^2.$$

Using $\mathbb{E} [VU^\top] = I$, we obtain

$$\mathbb{E} \left[V \left(\frac{h(\theta + \delta U) - h(\theta - \delta U)}{2\delta} \right) \right] = \mathbb{E} [VU^\top \nabla h(\theta) + \delta \phi(\theta, \delta, U)V]$$

$$= \nabla h(\theta) + \delta \hat{\phi}(\theta, \delta),$$

where $\hat{\phi}(\theta, \delta)$ satisfies $\|\hat{\phi}(\theta, \delta)\| \leq \frac{L}{2} \mathbb{E}[\|V\| \|U\|^2]$.

Applying the above expression to $F(\cdot, \psi)$ and using (3.23), we have

$$\mathbb{E} [\hat{\nabla} f(\theta)] = \nabla F(\theta, \psi) + \delta \hat{\phi}(\theta, \delta) \text{ a.s.},$$

where $\hat{\phi}(\theta, \delta)$ satisfies $\|\hat{\phi}(\theta, \delta)\| \leq \frac{L}{2} \mathbb{E}[\|V\| \|U\|^2]$.

A3.4 together with dominated convergence theorem leads to $E[\nabla F(\theta, \psi)] = \nabla f(\theta)$. Using this fact, we obtain

$$\begin{aligned} & \left\| \mathbb{E} [\hat{\nabla} f(\theta)] - \nabla f(\theta) \right\| \\ &= \left\| \mathbb{E} \left[V \left(\frac{f(\theta + \delta U) - f(\theta - \delta U)}{2\delta} \right) - V U^\top \nabla f(\theta) \right] \right\| \\ &\leq \delta \|\mathbb{E}[V \phi(\theta, \delta, U)]\| \\ &\leq \frac{\delta L}{2} \mathbb{E}[\|V\| \|U\|^2], \end{aligned}$$

and the claim for the bias follows by setting $C_1 = \frac{L}{2} \mathbb{E}[\|V\| \|U\|^2]$.

We now bound $\mathbb{E} \left[\|\hat{\nabla} f(\theta)\|^2 \right]$ as follows:

$$\begin{aligned} \mathbb{E} \|\hat{\nabla} f(\theta)\|^2 &\stackrel{\text{why?}}{=} \mathbb{E} \left\| V \left(\delta \phi(\theta, \delta, U) + U^\top \nabla f(\theta) \right) \right\|^2 \\ &\leq \mathbb{E} \left[\left(\|V U^\top \nabla f(\theta)\| + \frac{\delta L}{2} \|V\| \|U\|^2 \right)^2 \right] \\ &\leq 2\mathbb{E} \left[\|V U^\top \nabla f(\theta)\|^2 \right] + \frac{\delta^2 L^2}{2} \mathbb{E} \left[\|V\|^2 \|U\|^4 \right], \end{aligned}$$

and the claim for the variance follows by setting

$$C_2 = 2B_1^2 + \frac{L^2}{2} \mathbb{E} \left[\|V\|^2 \|U\|^4 \right] \text{ with } B_1 = \sup_{\theta} \|\nabla f(\theta)\|. \quad \square$$

$$\begin{aligned}
\mathbb{E} \left\| \widehat{\nabla} f(\theta) \right\|^2 &\stackrel{\text{Wu}}{=} \mathbb{E} \left\| V \left(\delta \phi(\theta, \delta, U) + U^\top \nabla f(\theta) \right) \right\|^2 \quad \text{--- } (\text{---}) \\
&\leq \mathbb{E} \left[\left(\|V U^\top \nabla f(\theta)\| + \frac{\delta L}{2} \|V\| \|U\|^2 \right)^2 \right] \\
&\leq 2\mathbb{E} \left[\|V U^\top \nabla f(\theta)\|^2 \right] + \frac{\delta^2 L^2}{2} \mathbb{E} \left[\|V\|^2 \|U\|^4 \right],
\end{aligned}$$

$$\mathbb{E} \left[\widehat{\nabla} f(\theta) \right] = \nabla F(\theta, \psi) + \delta \widehat{\phi}(\theta, \delta)$$

$$\mathbb{E} \left[V \left(\frac{h(\theta + \delta U) - h(\theta - \delta U)}{2\delta} \right) \right] = \mathbb{E} \left[V U^\top \nabla h(\theta) + \delta \phi(\theta, \delta, U) V \right] \quad \text{--- } (\text{---})$$

A similar claim with improved variance holds without assuming convexity.