
CHAPTER 5. ZERO-ORDER OPTIMIZATION

Yilong Song

October 17, 2023

1 Application Scenarios

1.1 Adversarial example generation for deep neural network

Let θ be a given DNN, from which we have complete access to the corresponding output of our input. However, the internal states/configurations and the operating mechanism of DL systems are not revealed to us. (Which calls a black-box)

The adversarial example generation problem can be described as:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & f(\mathbf{z} + \mathbf{x}; \theta) + \lambda g(\mathbf{x}) \\ \text{s.t.} \quad & \|\mathbf{x}\| \leq \epsilon, \end{aligned}$$

where \mathbf{z} denotes the former legitimate example and $\mathbf{z} + \mathbf{x}$ denotes an adversarial example, with the perturbation \mathbf{x} , f denotes the attack loss function, and $\lambda g(\mathbf{x})$ is a regularization function. Subject to the black-box θ , the accurate expression of f is not accessible, where zeroth-order methods apply well. [1]

1.2 Online sensor management

The problem could be described as:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & \frac{1}{T} \sum_{t=1}^T \left[-\log \det \left(\sum_{i=1}^d x_i \mathbf{a}_{i,t} \mathbf{a}_{i,t}^\top \right) \right] \\ \text{s.t.} \quad & \mathbf{1}^\top \mathbf{x} = m_0, \quad \mathbf{0} \leq \mathbf{x} \leq \mathbf{1}, \end{aligned}$$

where $\mathbf{x} \in \mathbb{R}^d$ is the optimization variable, d is the number of sensors, $\mathbf{a}_{i,t}$ is the observation coefficient of sensor i at time t , and m_0 is the number of selected sensors.

To calculate the gradient, we need the inverses of large matrices, which is costly.

2 Problem formulation

We list several situations where the first-order methods lose efficacy:

- Gradients are costly to compute.
- A black-box which keeps the accurate function expression invisible.
- The observed function value contains noise, which makes gradient not representative.

3 Standard Finite Differences

3.1 Gradient Estimator

Review: The partial derivative of f can be expressed as:

$$\frac{\partial f(\mathbf{x}_0)}{\partial x_i} = \lim_{t \rightarrow 0} \frac{f(\mathbf{x}_0 + t\mathbf{e}_i) - f(\mathbf{x}_0)}{t}, \quad i = 1, 2, \dots, d. \quad (1)$$

and the gradient is:

$$\nabla f(\mathbf{x}_0) = \left(\frac{\partial f(\mathbf{x}_0)}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x}_0)}{\partial x_d} \right)^\top = \sum_{i=1}^d \frac{\partial f(\mathbf{x}_0)}{\partial x_i} \mathbf{e}_i \quad (2)$$

However, the operation of limit is only theoretically feasible, for practice, an excessively small t for numerical computation brings a large scale of error. It is workable that we use a finite t instead of the limit of t to estimate gradient.

So we have a preliminary gradient estimator:

$$\hat{\nabla} f(\mathbf{x}_0)_t = \sum_{i=1}^d \frac{f(\mathbf{x}_0 + t\mathbf{e}_i) - f(\mathbf{x}_0)}{t} \mathbf{e}_i \quad (3)$$

In order to obtain an estimator of the gradient at one point, we need $d + 1$ samples of function value, namely $\mathcal{X} = \{\mathbf{x}_0, \mathbf{x}_0 + t\mathbf{e}_1, \dots, \mathbf{x}_0 + t\mathbf{e}_d\}$. The estimator above is called **forward difference**. Sometimes we also use a "two-side" estimator, which we call **central difference**:

$$\hat{\nabla} f(\mathbf{x}_0)_t = \sum_{i=1}^d \frac{f(\mathbf{x}_0 + t\mathbf{e}_i) - f(\mathbf{x}_0 - t\mathbf{e}_i)}{2t} \mathbf{e}_i \quad (4)$$

This method equals to a mean value on two opposite directions around \mathbf{x}_0 , namely $\mathcal{X} = \{\mathbf{x}_0 \pm t\mathbf{e}_1, \dots, \mathbf{x}_0 \pm t\mathbf{e}_d\}$, it requires $2d$ samples for one point, but usually this method is more accurate.

3.2 Error Analysis

Since a gradient without any limits could be inscrutable, we need some assumption for the gradient of f if we want to bound the error, although the gradient is not achievable in practice.

Assumption 3.1. (Lipschitz continuity of the gradients of f)

The function f is continuously differentiable, and the gradient of f is L -Lipschitz continuous for all $x \in \mathbb{R}^d$. Namely,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^d \quad (5)$$

Notation. This assumption ensures L -smooth of f , and it also guarantees that f is not too concave.

Theorem 3.2. Let $g(x)$ denote the forward finite difference approximation to the gradient $\nabla f(x)$, if $\nabla f(x)$ is L -Lipschitz continuous, then for all $x \in \mathbb{R}^d$,

$$\|g(x) - \nabla f(x)\| \leq \frac{\sqrt{d}Lt}{2} \quad (6)$$

Proof.

$$\begin{aligned} tg(x)^\top \mathbf{e}_i &= f(x + t\mathbf{e}_i) - f(x) \\ &= \int_0^t \frac{\partial}{\partial x_i} f(x + u\mathbf{e}_i) du = \int_0^t \mathbf{e}_i^\top \nabla f(x + u\mathbf{e}_i) du \end{aligned}$$

So,

$$\begin{aligned} t|\mathbf{e}_i^\top (g(x) - \nabla f(x))| &= \left| \int_0^t \mathbf{e}_i^\top \nabla f(x + u\mathbf{e}_i) du - t\mathbf{e}_i^\top \nabla f(x) \right| \\ &= \left| \int_0^t \mathbf{e}_i^\top (\nabla f(x + u\mathbf{e}_i) - \nabla f(x)) du \right| \\ &\leq \int_0^t \|\nabla f(x + u\mathbf{e}_i) - \nabla f(x)\| du \\ &\leq \int_0^t L u du = \frac{Lt^2}{2} \end{aligned}$$

So for every entry of $g(x) - \nabla f(x)$,

$$|(g(x) - \nabla f(x))_i| \leq \frac{Lt}{2}$$

then for the d -dimension vector,

$$\|g(x) - \nabla f(x)\| \leq \frac{\sqrt{d}Lt}{2}$$

□

Observation. It's obvious that the error would be smaller if we choose a smaller t , which will, however, bring in higher computational error.

3.3 Zeroth-order gradient descent with forward difference

With the finite difference estimators introduced above, we have a zeroth-order gradient descent (ZO-GD) method as follows

$$\mathbf{g}(\mathbf{x}_k) = \sum_{i=1}^d \frac{f(\mathbf{x}_k + t\mathbf{e}_i) - f(\mathbf{x}_k)}{t} \mathbf{e}_i \quad (7)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \gamma \mathbf{g}(\mathbf{x}_k) \quad (8)$$

3.3.1 Convergence in the non-convex scenario

Theorem 3.3. Assume $f(\mathbf{x})$ to be L -smooth. If we set $\gamma = 1/L$, ZO-GD converges as follows

$$\frac{1}{K+1} \sum_{k=0}^K \|\nabla f(\mathbf{x}_k)\|^2 \leq \frac{2L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{K+1} + \frac{dL^2t^2}{4}. \quad (9)$$

Proof. Since $f(\mathbf{x})$ is L -smooth, we have

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 \\ &= f(\mathbf{x}_k) - \gamma \langle \nabla f(\mathbf{x}_k), \mathbf{g}(\mathbf{x}_k) \rangle + \frac{\gamma^2 L}{2} \|\mathbf{g}(\mathbf{x}_k)\|^2 \\ &= f(\mathbf{x}_k) - \frac{\gamma}{2} \|\nabla f(\mathbf{x}_k)\|^2 - \frac{\gamma}{2} (1 - L\gamma) \|\mathbf{g}(\mathbf{x}_k)\|^2 + \frac{\gamma}{2} \|\mathbf{g}(\mathbf{x}_k) - \nabla f(\mathbf{x}_k)\|^2 \\ &\leq f(\mathbf{x}_k) - \frac{\gamma}{2} \|\nabla f(\mathbf{x}_k)\|^2 + \frac{\gamma dL^2t^2}{8} \end{aligned} \quad (10)$$

where the last inequality holds due to Theorem 3.2 and $\gamma \leq 1/L$. Taking the average of the above inequality over $k = 0, 1, \dots, K$, we achieve

$$\frac{1}{K+1} \sum_{k=0}^K \|\nabla f(\mathbf{x}_k)\|^2 \leq \frac{2(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{\gamma(K+1)} + \frac{dL^2t^2}{4}. \quad (11)$$

If we let $\gamma = 1/L$, then we prove the result. \square

3.3.2 Convergence in the strongly-convex scenario

Theorem 3.4. Assume $f(\mathbf{x})$ to be L -smooth and μ -strongly convex. If we set $\gamma = 1/L$, ZO-GD converges as follows

$$f(\mathbf{x}_K) - f(\mathbf{x}^*) \leq (1 - \frac{\mu}{L})^K (f(\mathbf{x}_0) - f(\mathbf{x}^*)) + \frac{dL^2t^2}{8\mu} \quad (12)$$

Proof. Since $f(x)$ is μ -strongly convex, it holds from Theorem 3.8 of our notes “Ch0” that

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2\mu} \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|^2. \quad (13)$$

Let $\mathbf{y} = \mathbf{x}_k$ and $\mathbf{x} = \mathbf{x}^*$, we have

$$\|\nabla f(\mathbf{x}_k)\|^2 \geq 2\mu(f(\mathbf{x}^k) - f(\mathbf{x}^*)). \quad (14)$$

Substituting the above inequality into (10), we achieve

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq (1 - \mu\gamma)(f(\mathbf{x}_k) - f(\mathbf{x}^*)) + \frac{\gamma d L^2 t^2}{8} \quad (15)$$

Keep iterating the above inequality, we have

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq (1 - \mu\gamma)^{k+1}(f(\mathbf{x}_0) - f(\mathbf{x}^*)) + \frac{d L^2 t^2}{8\mu} \quad (16)$$

Setting $\gamma = 1/L$ achieves the final result. \square

4 Linear Interpolation

4.1 Gradient Estimator

Standard Finite Differences use the sample set $\mathcal{X} = \{x, x + t\mathbf{e}_1, \dots, x + t\mathbf{e}_d\}$. Although gradients are always expressed in terms of entries when we need to calculate them, it's not necessary that we should always estimate the gradient along the directions of d axes. In fact, we can use a group of d linearly independent vectors $\{\mathbf{u}_i\}_{i=1}^d$ in stead of standard unit vectors. So the sample set now becomes $\mathcal{X} = \{x, x + t\mathbf{u}_1, \dots, x + t\mathbf{u}_d\}$. It's notable that \mathbf{u}_i needn't be a unit vector, but since the scale of \mathbf{u}_i is arbitrary, how can we construct a proper gradient estimator? We need the tool of **linear interpolation**. [2]

As Taylor's formula shows,

$$f(x + \Delta x) = f(x) + \nabla f(x)^\top \cdot \Delta x + o(\|\Delta x\|) \quad (17)$$

We extract the linear part,

$$m(x + \Delta x) = f(x) + \nabla f(x)^\top \cdot \Delta x \quad (18)$$

It's a d -dimension hyperplane in $(d + 1)$ -dimension space, which we need $d + 1$ points to determine when $\nabla f(x)$ is unknown. Precisely, we have $d + 1$ points in the sample set above. Namely, let $m(x + t\mathbf{u}_i) = f(x + t\mathbf{u}_i)$. We still use $\mathbf{g}(x)$ to represent the undetermined coefficients, the equation could be expressed as

$$\begin{pmatrix} f(x + t\mathbf{u}_1) - f(x) \\ \vdots \\ f(x + t\mathbf{u}_d) - f(x) \end{pmatrix} = t \begin{pmatrix} \mathbf{u}_1^\top \\ \vdots \\ \mathbf{u}_d^\top \end{pmatrix} \mathbf{g}(x) \quad (19)$$

or we simplify it as

$$F_{\mathcal{X}} = tQ_{\mathcal{X}}g(x) \quad (20)$$

So

$$\begin{aligned} g(x) &= \frac{1}{t}Q_{\mathcal{X}}^{-1}F_{\mathcal{X}} \\ &\triangleq \frac{1}{t}(\mathbf{v}_1, \dots, \mathbf{v}_d) \begin{pmatrix} f(x + t\mathbf{u}_1) - f(x) \\ \vdots \\ f(x + t\mathbf{u}_d) - f(x) \end{pmatrix} \\ &= \sum_{i=1}^d \frac{f(x + t\mathbf{u}_i) - f(x)}{t} \mathbf{v}_i \end{aligned}$$

where $\{\mathbf{v}_i\}$ are a group of generated vectors of $\{\mathbf{u}_i\}$ and satisfy

$$\begin{pmatrix} \mathbf{u}_1^\top \\ \vdots \\ \mathbf{u}_n^\top \end{pmatrix} (\mathbf{v}_1, \dots, \mathbf{v}_d) = \mathbf{I}_d \quad (21)$$

4.2 Error Analysis

Notation. For some \mathbf{u}_i whose scale is excessively large, we tend to involve the scalar in t and proportionally shrink every \mathbf{u}_i . So without loss of generality, we define $\|\mathbf{u}_i\| \leq 1$ for all $i = 1, 2, \dots$.

Theorem 4.1. Let $g(x)$ denote the forward linear interpolation approximation to the gradient $\nabla f(x)$, where $\mathcal{X} = \{x + t\mathbf{u}_1, \dots, x + t\mathbf{u}_d\}$ is a set of interpolation points such that $\max_{1 \leq i \leq d} \|\mathbf{u}_i\| \leq 1$ and $Q_{\mathcal{X}}$ is non-singular. If $\nabla f(x)$ is L -Lipschitz continuous, then for all $x \in \mathbb{R}^d$,

$$\|g(x) - \nabla f(x)\| \leq \frac{\|Q_{\mathcal{X}}^{-1}\|_2 \sqrt{d} L t}{2} \quad (22)$$

Proof.

$$\begin{aligned} tg(x)^\top \mathbf{u}_i &= f(x + t\mathbf{u}_i) - f(x) \\ &= \int_0^t \mathbf{u}_i^\top \nabla f(x + v\mathbf{u}_i) dv \end{aligned}$$

So,

$$\begin{aligned}
& t|\mathbf{u}_i^\top (g(x) - \nabla f(x))| \\
&= \left| \int_0^t \mathbf{u}_i^\top \nabla f(x + v\mathbf{u}_i) dv - t\mathbf{u}_i^\top \nabla f(x) \right| \\
&= \left| \int_0^t \mathbf{u}_i^\top (\nabla f(x + v\mathbf{u}_i) - \nabla f(x)) dv \right| \\
&\leq \int_0^t \|\mathbf{u}_i\| \|\nabla f(x + v\mathbf{u}_i) - \nabla f(x)\| dv \\
&\leq \int_0^t L \|\mathbf{u}_i\|^2 v dv = \frac{\|\mathbf{u}_i\|^2 Lt^2}{2} \leq \frac{Lt^2}{2}
\end{aligned}$$

which implies,

$$\|Q_{\mathcal{X}}(g(x) - \nabla f(x))\| \leq \frac{\sqrt{d}Lt}{2}$$

so,

$$\|g(x) - \nabla f(x)\| \leq \frac{\|Q_{\mathcal{X}}^{-1}\|_2 \sqrt{d}Lt}{2}$$

□

The result shows that large $\|Q_{\mathcal{X}}^{-1}\|$ can cause large deviation of $\mathbf{g}(x)$ from $\nabla f(x)$. Thus it is desirable to select \mathcal{X} in such a way that the condition number of $Q_{\mathcal{X}}^{-1}$ is small. It's the best when $Q_{\mathcal{X}}$ is orthonormal.

4.3 Central Difference

To make an analogue of central difference in standard finite difference, we use the sample set $\mathcal{X} = \{\mathbf{x}_0 \pm t\mathbf{u}_1, \dots, \mathbf{x}_0 \pm t\mathbf{u}_d\}$, and similarly let

$$\mathbf{g}(x) = \sum_{i=1}^d \frac{f(x + t\mathbf{u}_i) - f(x - t\mathbf{u}_i)}{2t} \mathbf{v}_i \quad (23)$$

Here we use $2d$ points to determine a hyperplane. It's actually the average of two hyperplanes.

5 Gaussian Smoothness

5.1 Random Estimator

Up to now, we always need $d + 1$ samples to estimate the gradient at one point, but if we choose just one direction for each iteration instead of d directions? We can certainly iterate through d selected directions during d iterations, but there is a more efficient method: random method.

To make the estimator unbiased, we need a central symmetrical variable \mathbf{u} (e.g. Gauss distribution or uniform distribution on a sphere), and along this direction, we preliminarily

define the gradient estimator as

$$\mathbf{g}(x) = \frac{f(x + t\mathbf{u}) - f(x)}{t} \mathbf{u} \quad (24)$$

It's notable that the direction of this estimator depends completely on the arbitrary \mathbf{u} , so it is far from a good estimator for the true gradient. However, when \mathbf{u} points at a bad direction, the factor before it may be small or negative. So it could hopefully be a proper estimator in terms of expectation.

5.2 Smooth Function

Sometimes the reason why we can't get the gradient is the noise, that is to say, the properties of f could be somewhat bad. But with the help of random method, we could construct a smooth version of f as

$$\hat{f}(x)_t \triangleq \mathbf{E}_{\mathbf{v}} f(x + t\mathbf{v}) \quad (25)$$

where t is called **smooth parameter**, and \mathbf{v} is another central symmetrical variable. This smooth function could be understood as the mean of the function value around $f(x)$. Usually the smooth version of f owes better properties than f .

5.3 Bridge Lemma for Gaussian Smoothness

Now we specify $\mathbf{u} \sim N(0, \mathbf{I})$, $\mathbf{v} \sim N(0, \mathbf{I})$, the following lemma bridges the gap between random estimator and smooth function, which ensures the feasibility of Gaussian smoothness.

Lemma 5.1. (bridge lemma I) If $\mathbf{u} \in \mathbb{R}^d$ and $\mathbf{u} \sim N(0, \mathbf{I})$, for $f \in \mathbb{C}^1$, a certain x , and a certain t ,

$$\mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x})}{t} \cdot \mathbf{u} = \nabla \mathbf{E}_{\mathbf{u}} f(\mathbf{x} + t\mathbf{u}) \quad (26)$$

Proof. Because

$$\mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x})}{t} \cdot \mathbf{u} = \begin{pmatrix} \mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x})}{t} \cdot u_1 \\ \vdots \\ \mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x})}{t} \cdot u_d \end{pmatrix},$$

for every entry, we have

$$\begin{aligned}
& \mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x})}{t} \cdot u_1 = \mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x} + t\mathbf{u})}{t} \cdot u_1 \\
&= \underbrace{\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty}}_d \frac{f(\mathbf{x} + t\mathbf{u})}{t} \cdot \frac{u_1}{(\sqrt{2\pi})^d} e^{-\frac{u_1^2 + \dots + u_d^2}{2}} du_1 \dots du_d \\
&= \underbrace{\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty}}_{d-1} \frac{1}{(\sqrt{2\pi})^{d-1}} e^{-\frac{u_2^2 + \dots + u_d^2}{2}} du_2 \dots du_d \int_{-\infty}^{+\infty} \frac{f(\mathbf{x} + t\mathbf{u})}{t} \cdot \frac{u_1}{\sqrt{2\pi}} e^{-\frac{u_1^2}{2}} du_1 \\
&= \underbrace{\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty}}_{d-1} \frac{1}{(\sqrt{2\pi})^{d-1}} e^{-\frac{u_2^2 + \dots + u_d^2}{2}} du_2 \dots du_d \left(- \int_{-\infty}^{+\infty} \frac{f(\mathbf{x} + t\mathbf{u})}{t} \cdot \frac{1}{\sqrt{2\pi}} de^{-\frac{u_1^2}{2}} \right) \\
&= \underbrace{\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty}}_{d-1} \frac{1}{(\sqrt{2\pi})^{d-1}} e^{-\frac{u_2^2 + \dots + u_d^2}{2}} du_2 \dots du_d \cdot \\
&\quad \left(- \frac{f(\mathbf{x} + t\mathbf{u})}{t} \cdot \frac{e^{-\frac{u_1^2}{2}}}{\sqrt{2\pi}} \Big|_{-\infty}^{+\infty} + \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u_1^2}{2}} \cdot \frac{\partial f(\mathbf{x} + t\mathbf{u})}{t \partial u_1} \cdot \frac{1}{\sqrt{2\pi}} du_1 \right) \\
&= \underbrace{\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty}}_{d-1} \frac{1}{(\sqrt{2\pi})^{d-1}} e^{-\frac{u_2^2 + \dots + u_d^2}{2}} du_2 \dots du_d \cdot \left(\int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u_1^2}{2}} \cdot f'_1(\mathbf{x} + t\mathbf{u}) \cdot \frac{1}{\sqrt{2\pi}} du_1 \right) \\
&= \underbrace{\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty}}_d f'_1(\mathbf{x} + t\mathbf{u}) \cdot \frac{1}{(\sqrt{2\pi})^d} e^{-\frac{u_1^2 + \dots + u_d^2}{2}} du_1 \dots du_d \\
&= \mathbf{E}_{\mathbf{u}} f'_1(\mathbf{x} + t\mathbf{u})
\end{aligned}$$

Therefore,

$$\begin{pmatrix} \mathbf{E}_{\mathbf{u}} f'_1(\mathbf{x} + t\mathbf{u}) \\ \vdots \\ \mathbf{E}_{\mathbf{u}} f'_d(\mathbf{x} + t\mathbf{u}) \end{pmatrix} = \mathbf{E}_{\mathbf{u}} \nabla f(\mathbf{x} + t\mathbf{u}) = \nabla \mathbf{E}_{\mathbf{u}} f(\mathbf{x} + t\mathbf{u}) \quad (27)$$

□

Notation. $f \in \mathbb{C}^1$ is just a sufficient condition, the lemma holds among a wide range of functions, thanks to the good properties of normal distribution.

The lemma demonstrates that the expectation of our gradient estimator is precisely the gradient of the smooth function, which proves theoretically that our estimator is somehow an unbiased estimation of gradient. So the estimator could be denoted by

$$\mathbf{g}(x) = \frac{f(x + t\mathbf{u}) - f(x)}{t} \mathbf{u}, \quad \mathbf{u} \sim N(0, \mathbf{I}) \quad (28)$$

5.4 Two-point Estimator vs Single-point Estimator

From the proof above, we notice that $\mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x}+t\mathbf{u})-f(\mathbf{x})}{t} \cdot \mathbf{u} = \mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x}+t\mathbf{u})}{t} \cdot \mathbf{u}$, since the expectation of the second term is zero. We name the former form **two-point estimator**, while the latter **single-point estimator**. These two forms of estimator are the same in terms of expectation, but different in practice. Obviously, single-point estimator will explode to ∞ as $t \rightarrow 0$, and it owes naturally a larger variance. The convergence analysis of single-point estimator is also different from the two-point estimator.

5.5 Error Analysis in Terms of Probability

The difference between our estimator and the true gradient can be divided into two parts:

$$\|\mathbf{g}(x) - \nabla f(x)\| \leq \|\nabla \hat{f}_t(x) - \nabla f(x)\| + \|\mathbf{g}(x) - \nabla \hat{f}_t(x)\| \quad (29)$$

The first term is the difference between the true function and the smooth function, which can be bounded as:

Lemma 5.2. if $\nabla f(x)$ is L -Lipschitz continuous, then for all $x \in \mathbb{R}^d$,

$$\|\nabla \hat{f}_t(x) - \nabla f(x)\| \leq \sqrt{dLt} \quad (30)$$

The second term is caused by the estimator. In fact, the covariance of $\mathbf{g}(x)$ could be enormously huge if we use only one sample. To reduce it, we can take N samples and calculate their mean.

$$\mathbf{g}(x) = \frac{1}{N} \sum_{i=1}^N \frac{f(x+t\mathbf{u}_i) - f(x)}{t} \mathbf{u}_i, \quad \mathbf{u}_i \sim N(0, \mathbf{I}) \quad (31)$$

If we want to let the second term of error less than a determined r with probability $1 - \delta$, we need a large enough N . Considering both terms, the following theorem gives a lower bound of N .

Theorem 5.3. Suppose that $\nabla f(x)$ is L -Lipschitz continuous, if

$$N \geq \frac{3d}{\delta r^2} \left(3\|\nabla f(x)\|^2 + \frac{L^2 t^2}{4}(d+2)(d+4) \right) \quad (32)$$

then for all $x \in \mathbb{R}^d$ and $r > 0$,

$$\|\mathbf{g}(x) - \nabla f(x)\| \leq \sqrt{dLt} + r \quad (33)$$

with probability at least $1 - \delta$.

5.6 Error Analysis in Terms of Expectation

Lemma 5.4. Suppose f is L -smooth, and let $\mathbf{u} \sim N(\mathbf{0}, \mathbf{I})$. Then,

$$\mathbf{E}_{\mathbf{u}}[\|\mathbf{g}(x)\|^2] \leq 2(d+2)\|\nabla f(x)\|^2 + \frac{t^2 L^2}{2}d(d+2)(d+4) \quad (34)$$

Proof. We have

$$\begin{aligned} \mathbf{E}_{\mathbf{u}}\|\mathbf{g}(x)\|^2 &= \mathbf{E}_{\mathbf{u}}\left\|\frac{f(x+t\mathbf{u})-f(x)}{t} \cdot \mathbf{u}\right\|^2 \\ &= \frac{1}{t}\mathbf{E}_{\mathbf{u}}[\|f(x+t\mathbf{u})-f(x)\|^2\|\mathbf{u}\|^2], \end{aligned}$$

considering $(a+b)^2 \leq 2a^2 + 2b^2$, we have

$$\mathbf{E}_{\mathbf{u}}\|\mathbf{g}(x)\|^2 \leq \frac{2}{t^2}\mathbf{E}_{\mathbf{u}}[|f(x+t\mathbf{u})-f(x)-\langle\nabla f(x), t\mathbf{u}\rangle|^2\|\mathbf{u}\|^2] + 2\mathbf{E}_{\mathbf{u}}[|\nabla f(x), \mathbf{u}|^2\|\mathbf{u}\|^2].$$

First we consider the second term, note that

$$\mathbf{E}_{\mathbf{u}}[|\langle\nabla f(x), \mathbf{u}\rangle\|\mathbf{u}\|^2] = (\nabla f(x))^\top \mathbf{E}_{\mathbf{u}}\|\mathbf{u}\|^2 \mathbf{u} \mathbf{u}^\top \nabla f(x). \quad (35)$$

When $i \neq j$, $\mathbf{E}_{\mathbf{u}}\|\mathbf{u}\|^2 u_i u_j = 0$, for the independence and central symmetry of every variable. When $i = j$,

$$\mathbf{E}_{\mathbf{u}}\|\mathbf{u}\|u_i^2 = \mathbf{E}_{\mathbf{u}}(u_1^2 + \dots + u_d^2)u_i^2 = \mathbf{E}_{\mathbf{u}}u_1^4 + (d-1)(\mathbf{E}_{\mathbf{u}}u_1^2)^2 = 3 + (d-1) = d+2. \quad (36)$$

Therefore,

$$\mathbf{E}_{\mathbf{u}}[\|\mathbf{u}\|^2 \mathbf{u} \mathbf{u}^\top] = (d+2)\mathbf{I}. \quad (37)$$

$$2\mathbf{E}_{\mathbf{u}}[|\langle\nabla f(x), \mathbf{u}\rangle|^2\|\mathbf{u}\|^2] = 2(d+2)\|\nabla f(x)\|^2 \quad (38)$$

Next we bound the first term: By Newton-Leibniz theorem,

$$f(x+t\mathbf{u})-f(x) = \int_0^t \frac{d}{ds}[f(x+s\mathbf{u})-f(x)]ds = \int_0^t \langle\nabla f(x+s\mathbf{u}), \mathbf{u}\rangle ds,$$

and thus,

$$\begin{aligned} |f(x+t\mathbf{u})-f(x)-\langle\nabla f(x), t\mathbf{u}\rangle| &= \int_0^t \langle\nabla f(x+s\mathbf{u})-\nabla f(x), \mathbf{u}\rangle ds \\ &\leq \int_0^t \|\nabla f(x+s\mathbf{u})-\nabla f(x)\|\|\mathbf{u}\|ds \\ &\leq \int_0^t Ls\|\mathbf{u}\|^2 ds = \frac{Lt^2}{2}\|\mathbf{u}\|^2 \end{aligned}$$

We then get:

$$\frac{2}{t^2}\mathbf{E}_{\mathbf{u}}[|f(x+t\mathbf{u})-f(x)-\langle\nabla f(x), t\mathbf{u}\rangle|^2\|\mathbf{u}\|^2] \leq \frac{2}{t^2}\mathbf{E}_{\mathbf{u}}\left[\frac{L^2 t^4}{4}\|\mathbf{u}^6\|\right] = \frac{L^2 t^2}{2}d(d+2)(d+4),$$

where $\mathbf{E}_{\mathbf{u}}\|\mathbf{u}\|^6$ is the third moment of χ^2 distribution.

□

6 Sphere Smoothness

6.1 Bridge Lemma for Sphere Smoothness

Now we consider the uniform distribution on a d -dimension spherical surface, denoted by $\mathbf{u} \sim U(\mathbb{S}^{d-1}(0, 1))$. Namely, \mathbf{u} is a unit vector, and the direction is central symmetrical. When we use it to construct our gradient estimation, the according generated \mathbf{v} is no longer the same as \mathbf{u} , but obeys uniform distribution in a unit ball, denoted by $\mathbf{v} \sim U(\mathbb{B}^d(0, 1))$, coupled with a coefficient d . The following lemma proves it. (*We need the knowledge of integral on manifolds to accomplish our proof.) [3]

Lemma 6.1. (bridge lemma II) If $\mathbf{u} \in \mathbb{R}^d$ and $\mathbf{u} \sim U(\mathbb{S}^{d-1}(0, 1))$, $\mathbf{v} \in \mathbb{R}^d$ and $\mathbf{v} \sim U(\mathbb{B}^d(0, 1))$, for $f \in \mathbb{C}^1$, a certain x , and a certain t ,

$$\mathbf{E}_{\mathbf{u}} \frac{f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x})}{t} \cdot \mathbf{u} = \frac{1}{d} \nabla \mathbf{E}_{\mathbf{v}} f(\mathbf{x} + t\mathbf{v}) \quad (39)$$

Proof. *

If $d = 1$, then the fundamental theorem of calculus implies,

$$\frac{d}{dx} \int_{-t}^t f(x + v) dv = f(x + t) - f(x - t) \quad (40)$$

The d -dimensional generalization, which follows from Stoke's theorem, is,

$$\nabla \int_{t\mathbb{B}} f(\mathbf{x} + \mathbf{v}) d\mathbf{v} = \int_{t\mathbb{S}} f(\mathbf{x} + \mathbf{u}) \frac{\mathbf{u}}{\|\mathbf{u}\|} d\mathbf{u} \quad (41)$$

By definition,

$$\mathbf{E} f(\mathbf{x} + t\mathbf{v}) = \frac{\int_{t\mathbb{B}} f(\mathbf{x} + \mathbf{v}) d\mathbf{v}}{\text{vol}_d(t\mathbb{B})} \quad (42)$$

Similarly,

$$\mathbf{E} f(\mathbf{x} + t\mathbf{u}) \cdot \mathbf{u} = \frac{\int_{t\mathbb{S}} f(\mathbf{x} + \mathbf{u}) \cdot \frac{\mathbf{u}}{\|\mathbf{u}\|} d\mathbf{u}}{\text{vol}_{d-1}(t\mathbb{S})} \quad (43)$$

Noticed that $\text{vol}_d(t\mathbb{B})$ is $\frac{t}{d}$ times of $\text{vol}_{d-1}(t\mathbb{S})$, then the lemma follows. □

So an unbiased estimator of sphere smoothness should be:

$$\mathbf{g}(x) = \frac{f(\mathbf{x} + t\mathbf{u}) - f(\mathbf{x})}{t} \cdot d \cdot \mathbf{u}, \quad \mathbf{u} \sim U(\mathbb{S}^{d-1}(0, 1)) \quad (44)$$

where d is an extra coefficient.

6.2 Error Analysis in Terms of Probability

To make an analogue, we also use a minibatch of samples as:

$$\mathbf{g}(x) = \frac{d}{N} \sum_{i=1}^N \frac{f(x + t\mathbf{u}) - f(x)}{t} \cdot \mathbf{u} \quad \mathbf{u} \sim U(\mathbb{S}^{d-1}(0, 1)) \quad (45)$$

And the similar theorem is as follows

Theorem 6.2. Suppose that $\nabla f(x)$ is L -Lipschitz continuous, if

$$N \geq \left[\frac{6d^2}{r^2} \left(\frac{\|\nabla f(x)\|^2}{d} + \frac{L^2 t^2}{4} \right) + \frac{2d}{3r} (2\|\nabla f(x)\| + Lt) \right] \log \frac{d+1}{\delta} \quad (46)$$

then for all $x \in \mathbb{R}^d$ and $r > 0$,

$$\|\mathbf{g}(x) - \nabla f(x)\| \leq Lt + r \quad (47)$$

with probability at least $1 - \delta$.

6.3 Error Analysis in Terms of Expectation

Lemma 6.3. Suppose f is L -smooth, and let $\mathbf{u} \sim U(\mathbb{S}^{d-1}(0, 1))$. Then,

$$\mathbf{E}_{\mathbf{u}}[\|\mathbf{g}(x)\|^2] \leq 2d\|\nabla f(x)\|^2 + \frac{t^2 L^2 d^2}{2} \quad (48)$$

The proof is similar to the parallel lemma for Gaussian smoothness, but we need the fact that:

$$\mathbf{E}_{\mathbf{u}}[\mathbf{u}\mathbf{u}^\top] = \frac{1}{d}\mathbf{I}, \quad \mathbf{u} \sim U(\mathbb{S}^{d-1}(0, 1)), \quad (49)$$

which also demands the integral on a manifold.

6.4 Convergence Analysis in convex situation

Now we assume that $f(x)$ is convex and L -smooth, we analyze the convergence rate of Sphere Smoothness. With the assumption of convexity, the error between f and the smooth version of f has nothing to do with the number of dimension any more. [4]

Lemma 6.4. Suppose f is convex and L -smooth. Let $\mathbf{u} \sim N(\mathbf{0}, \mathbf{I})$, and \hat{f}_t denote the corresponding smooth version of f . Then \hat{f}_t is convex, L -smooth and satisfies

$$f(x) \leq \hat{f}_t(x) \leq f(x) + \frac{Lt^2}{2} \quad (50)$$

and

$$\|\nabla \hat{f}_t(x) - \nabla f(x)\| \leq Lt \quad (51)$$

Proof. The convexity of \hat{f}_t follows by noting that

$$\begin{aligned}
\hat{f}_t(\theta x_1 + (1 - \theta)x_2) &= \mathbf{E}_{\mathbf{v}}[f(\theta x_1 + (1 - \theta)x_2 + t\mathbf{v})] \\
&= \mathbf{E}_{\mathbf{v}}[f(\theta(x_1 + t\mathbf{v}) + (1 - \theta)(x_2 + t\mathbf{v}))] \\
&\leq \mathbf{E}_{\mathbf{v}}[\theta f(x_1 + t\mathbf{v}) + (1 - \theta)f(x_2 + t\mathbf{v})] \\
&= \theta \hat{f}_t(x_1) + (1 - \theta)\hat{f}_t(x_2)
\end{aligned}$$

for any $\theta \in [0, 1]$ and any x_1, x_2 . To show the L -smoothness of \hat{f}_t , let $x_1, x_2 \in \mathbb{R}^d$ be arbitrary, and we have

$$\begin{aligned}
\|\nabla \hat{f}_t(x_1) - \nabla \hat{f}_t(x_2)\| &= \|\nabla \mathbf{E}_{\mathbf{v}}[f(x_1 + t\mathbf{v})] - \nabla \mathbf{E}_{\mathbf{v}}[f(x_2 + t\mathbf{v})]\| \\
&= \|\mathbf{E}_{\mathbf{v}}[\nabla f(x_1 + t\mathbf{v}) - \nabla f(x_2 + t\mathbf{v})]\| \\
&\leq \mathbf{E}_{\mathbf{v}}[\|\nabla f(x_1 + t\mathbf{v}) - \nabla f(x_2 + t\mathbf{v})\|] \\
&\leq \mathbf{E}_{\mathbf{v}}[L\|x_1 - x_2\|] = L\|x_1 - x_2\|.
\end{aligned}$$

Now by the convexity and smoothness of f , we have

$$f(x) + \langle \nabla f(x), t\mathbf{v} \rangle \leq f(x + t\mathbf{v}) \leq f(x) + \langle \nabla f(x), t\mathbf{v} \rangle + \frac{L}{2}\|t\mathbf{v}\|^2 \quad (52)$$

Now we take the expectation with respect to $\mathbf{v} \sim N(\mathbf{0}, \mathbf{I})$, we have $\mathbf{E}_{\mathbf{v}}[\langle \nabla f(x), t\mathbf{v} \rangle] = 0$. Therefore,

$$f(x) \leq \mathbf{E}_{\mathbf{v}}[f(x + t\mathbf{v})] \leq f(x) + \frac{Lt^2}{2}\mathbf{E}_{\mathbf{v}}[\|\mathbf{v}\|^2], \quad (53)$$

which gives the first inequality. Now we regarding $\nabla \hat{f}_t$, we have

$$\begin{aligned}
\|\nabla \hat{f}_t(x) - \nabla f(x)\| &= \|\nabla_x \mathbf{E}_{\mathbf{v}}[f(x + t\mathbf{v}) - f(x)]\| \\
&= \|\mathbf{E}_{\mathbf{v}}[\nabla_x f(x + t\mathbf{v}) - \nabla_x f(x)]\| \\
&\leq \mathbf{E}_{\mathbf{v}}[\|\nabla f(x + t\mathbf{v}) - \nabla f(x)\|] \\
&\leq \mathbf{E}_{\mathbf{v}}[L\|t\mathbf{v}\|] \leq Lt,
\end{aligned}$$

□

We now make an iteration as follows, note that now t_k is not a constant, but a convergent series.

$$\mathbf{g}(x) = \frac{d}{t_k}(f(x + t_k \mathbf{u}) - f(x))\mathbf{u}, \quad \mathbf{u} \sim U(\mathbb{S}^{d-1}(0, 1)) \quad (54)$$

$$x_{k+1} = x_k - \gamma \mathbf{g}(x) \quad (55)$$

We have the theorem for convergence rate for Sphere Smoothness as follows:

Theorem 6.5. Suppose f is convex and L -smooth, and has a minimizer $x \in \mathbb{R}^d$. Let $\gamma = c/(dL)$ for some $c \in (0, 1)$, and let t_k be a positive sequence of smoothing parameter such that $\sum_{k=0}^K t_k^2 = R^2 \leq +\infty$. Then the zeroth-order optimization iteration achieves

$$\frac{1}{K+1} \sum_{k=0}^K \mathbf{E}[f(x_k) - f(x^*)] \leq \frac{dL\|x_0 - x^*\|^2}{c(1-c)(K+1)} + \frac{L}{2(1-c)} \left(1 + \frac{cd}{4}\right) \frac{\sum_{k=0}^K t_k^2}{K+1} \quad (56)$$

Proof.

$$\begin{aligned} \|x_{k+1} - x^*\|^2 &= \|x_k - x^* - \gamma \mathbf{g}(x)\|^2 \\ &= \|x_k - x^*\|^2 - 2\gamma \langle x_k - x^*, \mathbf{g}(x) \rangle + \gamma^2 \|\mathbf{g}(x)\|^2. \end{aligned}$$

Using the lemma for $\|\mathbf{g}(x)\|^2$ and the lemma for $\nabla \hat{f}_t$, we get

$$\mathbf{E} \langle x_k - x^*, \mathbf{g}(x) \rangle = \langle x_k - x^*, \nabla \hat{f}_{t_k}(x_k) \rangle, \quad (57)$$

$$\mathbf{E} \|\mathbf{g}(x)\|^2 \leq 2d \|\nabla f(x_k)\|^2 + \frac{t_k^2 L^2 d^2}{2}. \quad (58)$$

and consequently,

$$\mathbf{E} \|x_{k+1} - x^*\|^2 \leq \|x_k - x^*\|^2 - 2\gamma \langle x_k - x^*, \nabla \hat{f}_{t_k}(x_k) \rangle + 2\gamma^2 d \|\nabla f(x_k)\|^2 + \frac{\gamma^2 t_k^2 L^2 d^2}{2}. \quad (59)$$

Since \hat{f}_{t_k} is convex, we see that

$$\hat{f}_{t_k}(x^*) - \hat{f}_{t_k}(x_k) \geq \langle \nabla \hat{f}_{t_k}(x_k), x^* - x_k \rangle, \quad (60)$$

and coupled with the lemma, we have

$$-\langle \nabla \hat{f}_{t_k}(x_k), x^* - x_k \rangle \leq \hat{f}_{t_k}(x^*) - \hat{f}_{t_k}(x_k) \leq f(x^*) - f(x_k) + \frac{Lt_k^2}{2}. \quad (61)$$

Moreover, since f is L -smooth and $\nabla f(x^*) = 0$, we have

$$\|\nabla f(x_k)\|^2 = \|\nabla f(x_k) - \nabla f(x^*)\|^2 \leq 2L(f(x_k) - f(x^*)). \quad (62)$$

Summarizing these results, we get

$$\mathbf{E} \|x_{k+1} - x^*\|^2 \leq \|x_k - x^*\|^2 - 2\gamma(1 - 2\gamma dL)(f(x_k) - f(x^*)) + \gamma L t_k^2 + \frac{\gamma^2 t_k^2 L^2 d^2}{2}, \quad (63)$$

and by taking the total expectation, we can get

$$2\gamma(1 - 2\gamma dL) \mathbf{E}[f(x_k) - f(x^*)] \leq \mathbf{E} \|x_k - x^*\|^2 - \mathbf{E} \|x_{k+1} - x^*\|^2 + \gamma L t_k^2 + \frac{\gamma^2 t_k^2 L^2 d^2}{2}. \quad (64)$$

Now we take the telescoping sum and get

$$2\gamma(1 - 2\gamma dL) \sum_{k=0}^K \mathbf{E}[f(x_k) - f(x^*)] \leq \|x_0 - x^*\|^2 + \gamma L \left(1 + \frac{\gamma L d^2}{2}\right) \sum_{k=0}^K t_k^2. \quad (65)$$

By taking $\gamma = c/(2dL)$ for some $c \in (0, 1)$, we get the result:

$$\frac{1}{K+1} \sum_{k=0}^K \mathbf{E}[f(x_k) - f(x^*)] \leq \frac{dL\|x_0 - x^*\|^2}{c(1-c)(K+1)} + \frac{L}{2(1-c)} \left(1 + \frac{cd}{4}\right) \frac{\sum_{k=0}^K t_k^2}{K+1} \quad (66)$$

$$= \frac{d}{K+1} \left(\frac{L\|x_0 - x^*\|^2}{c(1-c)} + \frac{R^2 L(c + 4/d)}{8(1-c)} \right) \quad (67)$$

□

The result shows the convergence rate of zeroth-order optimization method with sphere smoothness is $O(d/K)$, which is still concerned with the number of dimension d .

References

- [1] L. Sijia, C. Pin-Yu, K. Bhavya, Z. Gaoyuan, H. I. Alfred O., and V. Pramod K., “A primer on zeroth-order optimization in signal processing and machine learning - principals, recent advances, and applications,” *Nonconvex Optimization for Signal Processing and Machine Learning*, pp. 43–54.
- [2] B. Albert S., C. Liyuan, C. Krzysztof, and S. Katya, “A theoretical and empirical comparison of gradient approximations in derivative-free optimization,” *Found Comput Math* 22, pp. 507–560.
- [3] F. Abraham D., K. Adam Tauman, and M. H. Brendan, “Online convex optimization in the bandit setting: gradient descent without a gradient,” *Proceedings of the sixteenth and annual ACM-SIAM symposium on discrete algorithms*, pp. 385–394.
- [4] T. Yujie, “Introduction to zeroth-order optimization.”