

Clarke Differential

W

Clarke Differential

$$\left(\begin{array}{l} x_{t+1} \leftarrow x_t - g_t y_t \\ g_t \in \partial f(x_t) \end{array} \right)$$

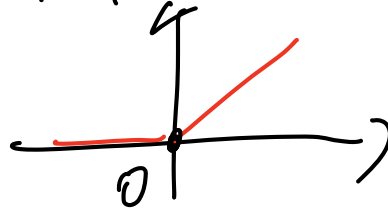
Definition: Given $f: \mathbb{R}^d \rightarrow \mathbb{R}$, for every x , the Clarke differential is defined as

$$\partial f(x) \triangleq \underline{\text{conv}} \left(\{s \in \mathbb{R}^d : \exists \{x_i\}_{i=1}^{\infty} \rightarrow x, \{\nabla f(x_i)\}_{i=1}^{\infty} \rightarrow s\} \right).$$

The elements in the subdifferential set are subgradients.

$$\text{conv}(S) = \left\{ v : v = \sum_{i=1}^n \lambda_i u_i, u_i \in S, \sum_{i=1}^n \lambda_i = 1, \lambda_i \geq 0 \right\}$$

example: ReLU



$$\{x_i\}, -1, -\frac{1}{2}, \dots \rightarrow 0$$

$$\nabla f(x_i) = 0$$

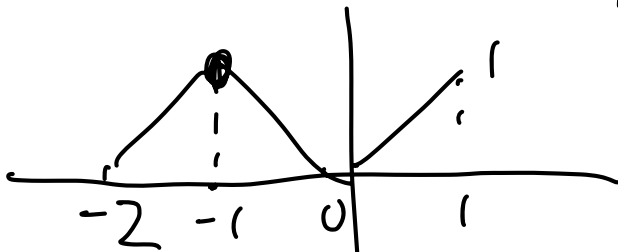
$$\{x_i\}, 1, \frac{1}{2}, \dots \rightarrow 0$$

$$\nabla f(x_i) = 1$$

$$\Rightarrow \partial f(0) = \left\{ \begin{array}{l} \lambda_1 \nabla f(x_1) \\ \lambda_2 \nabla f(x_2) \\ \lambda_1 + \lambda_2 = 1 \\ \lambda_1, \lambda_2 \geq 0 \end{array} \right\}$$

$$= [0, 1]$$

example



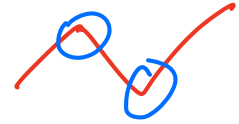
$$\{x_i\}: -2, -1.5, \dots \rightarrow \uparrow, \nabla f(x_i) = 1$$

$$\{x_i\}: 0, -0.5, \dots \rightarrow \uparrow, \nabla f(x_i) = -1$$

$$\Rightarrow \partial f(-1)$$

$$= [-1, 1],$$

When does Clarke differential exists



$$|f(x) - f(x')| \leq L \|x - x'\|, \quad x' \in \text{neighbor}(x)$$

Definition (Locally Lipschitz): $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is locally Lipschitz if $\forall x \in \mathbb{R}^d$, there exists a neighborhood S of x , such that f is Lipschitz in S .

① If f is locally Lip \Rightarrow Clarke differential exists
 \Rightarrow NN with ReLU exist
have Clarke differential everywhere

- If f is convex $\Rightarrow \partial f = \partial_s f$

- If f is differentiable $\Rightarrow \partial f = \{\nabla f\}$

~~*~~ satisfies chain rule

Positive Homogeneity

motiviert $\mathbb{R}eLU$

Definition: $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is positive homogeneous of degree L if $f(\alpha x) = \alpha^L f(x)$ for any $\alpha \geq 0$.

(1) $\mathbb{R}eLU: \sigma(\alpha z) = \alpha \cdot \sigma(z)$

(2) monomials: $\prod_{i=1}^d x_i^{p_i}, \sum_{i=1}^d p_i = L$

$$\begin{aligned} \prod_{i=1}^d (\alpha x_i)^{p_i} &= \alpha^{\sum p_i} \prod_{i=1}^d x_i^{p_i} \\ &= \alpha^L \prod_{i=1}^d x_i^{p_i} \end{aligned}$$

(3) Norm: $\|\alpha x\| = \alpha \cdot \|x\|$

Positive Homogeneity

(4) Multi-layer ReLU

$$f(x, W_1, \dots, W_{H+1}) = W_{H+1} \delta(W_H \delta(\dots \delta(W_1 x) \dots))$$

for one-layer, degree-1

$$f(x, W_1, \dots, \alpha W_H, \dots, W_{H+1}) = \alpha W_{H+1} \delta(W_H \dots \delta(W_1 x) \dots)$$

for all-layers

$$f(x, \alpha W_1, \dots, \alpha W_{H+1}) = \alpha^{H+1} f(x, W_1, \dots, W_{H+1})$$

\Rightarrow degree-(H+1)

Positive Homogeneity

$$W_n \in \mathcal{R}^{m \times n}$$

Fact: $\forall h = 1, \dots, H+1$

$$\left\langle W_n, \frac{\partial f(x, \dots, W_{H+1})}{\partial W_n} \right\rangle = f(x, W_1, \dots, W_{H+1})$$

pf: $A_n = \text{diag} (\sigma'(W_n) \sigma(\dots \sigma(W_1) \dots)) \in \mathcal{R}^{m \times n}$

$\sigma' = 0$ or $1 \Rightarrow$ (all or whether activation on or not)

$$f(x, W_1, \dots, W_{H+1}) = \underbrace{W_{H+1} A_H W_H \dots A_1 W_1 x}$$

$$\sigma(z) = \sigma'(z) \cdot z$$

$$\frac{\partial f(x, \dots, W_{H+1})}{\partial W_n} = (W_{H+1} A_H \dots W_{n+1} A_n)^T (A_{n+1} \dots W_1 x)$$

$$\left\langle W_n, \frac{\partial f}{\partial W_n} \right\rangle = f(x, \dots, W_{H+1})$$

Positive Homogeneity and Clark Differential

\Rightarrow Clarke differential exist

Lemma: Suppose $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is Locally Lipschitz and L -positively homogeneous. For any $x \in \mathbb{R}^d$ and $s \in \partial f(x)$, we have $\langle s, x \rangle = Lf(x)$.

Norm Preservation

$$f(X, W_1, W_2, W_3) = W_3 \sigma(W_2 \sigma(W_1 X))$$

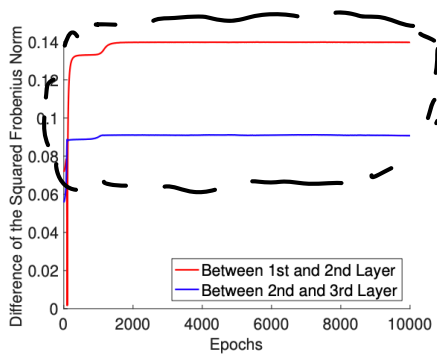
quadratic loss

$$\|W\|_F^2 = \sum_{i,j} W_{ij}^2$$

$$W_3 = 10$$

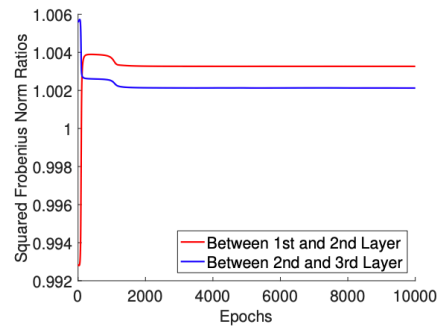
$$W_2 = 10$$

$$W_2(t+1) \leftarrow W_2(t) - \eta \frac{\partial L}{\partial W_2}$$



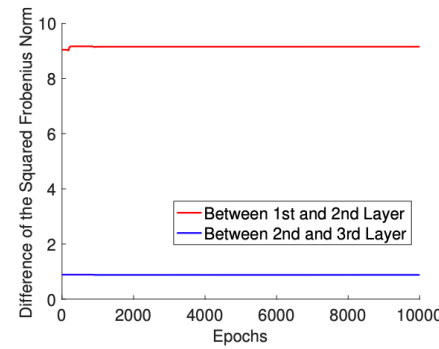
(a) Balanced initialization, squared norm differences.

$$\begin{aligned} & \|W_1\|_F^2 - \|W_2\|_F^2 \\ & \|W_2\|_F^2 - \|W_3\|_F^2 \end{aligned}$$

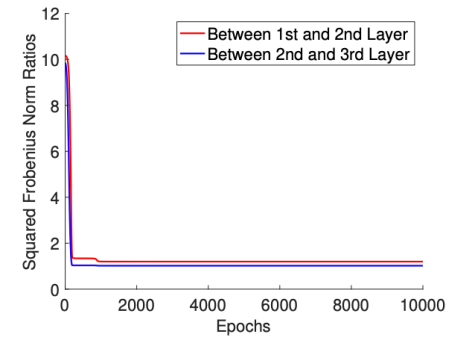


(b) Balanced initialization, squared norm ratios.

$$\frac{\|W_1\|_F^2}{\|W_2\|_F^2}$$



(c) Unbalanced Initialization, squared norm differences.



(d) Unbalanced initialization, squared norm ratios.

Gradient flow and gradient inclusion

Discrete-time dynamics can be complex. Let's use continuous-time dynamics to simplify:

$$\text{Gradient flow: } x_{t+1} = x_t - \eta \nabla f(x_t) \Rightarrow \frac{dx(t)}{dt} = -\nabla f(x(t))$$

$$\text{Gradient inclusion: } \frac{dx(t)}{dt} \in \partial f(x(t))$$

$$\frac{x_{t+\eta} - x_t}{\eta} = -\sigma f(x)$$

let $\eta \rightarrow 0$

Norm preservation by gradient inclusion

Algorithm regularization if $\|W_i(0)\|_F^2$ small for i

$\Rightarrow \|W_i(t)\|_F^2 \approx \|W_j(t)\|_F^2$

Theorem (Du, Hu, Lee '18) Suppose $\alpha > 0$, $f(x; (W_{H+1}, \dots, \alpha W_i, \dots, W_1)) = \alpha f(x, (W_{H+1}, \dots, W_1))$, i.e., predictions are 1-homogeneous in each layer. Then for every pair of layers $(i, j) \in [H+1] \times [H+1]$, the gradient inclusion maintains: for all $t \geq 0$,

$$\frac{1}{2} \|W_i(t)\|_F^2 - \frac{1}{2} \|W_i(0)\|_F^2 = \frac{1}{2} \|W_j(t)\|_F^2 - \frac{1}{2} \|W_j(0)\|_F^2.$$

$\Rightarrow \|W_i(t)\|_F^2 - \|W_j(t)\|_F^2 = \|W_i(0)\|_F^2 - \|W_j(0)\|_F^2$

Small in the init

Pf: (1) $\frac{dW_i(t)}{dt}$ formula

$$(2) \frac{1}{2} \|W_i(t)\|_F^2 - \frac{1}{2} \|W_i(0)\|_F^2 = \int_0^t \frac{d}{dt} \frac{1}{2} \|W_i(t)\|_F^2 dt$$

Optimization Methods for Deep Learning



Gradient descent for non-convex optimization

Descent Lemma: Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be twice differentiable, and $\|\nabla^2 f\|_2 \leq \beta$. Then setting the learning rate $\eta = 1/\beta$, and applying gradient descent, $x_{t+1} = x_t - \eta \nabla f(x_t)$, we have:

$$f(x_t) - f(x_{t+1}) \geq \frac{1}{2\beta} \|\nabla f(x_t)\|_2^2.$$

Pf: by Taylor expansion & mean-value theorem
$$f(x+\Delta) = f(x) + \Delta^T \nabla f(x) + \frac{1}{2} \Delta^T \nabla^2 f(\gamma) \Delta$$

(a.d. $f(x+\Delta) \approx f(x) + \Delta^T \nabla f(x) + \frac{1}{2} \|\Delta\|_2^2$) for some γ
 \Rightarrow optimize $\Delta \rightarrow -\nabla f(x)$

$$\Delta^T \nabla^2 f(\gamma) \Delta \leq \beta \|\Delta\|_2^2, \text{ choose } \Delta = -\eta \nabla f(x)$$
$$f(x_{t+1}) \leq f(x_t) - \eta \|\nabla f(x_t)\|_2^2 + \frac{1}{2} \beta \eta^2 \|\nabla f(x_t)\|_2^2 \leq f(x_t) - \frac{\eta}{2\beta} \|\nabla f(x_t)\|_2^2$$

Converging to stationary points

approximate
stationary point

Theorem: In $T = O\left(\frac{\beta}{\epsilon^2}\right)$ iterations, we have $\|\nabla f(x)\|_2 \leq \epsilon$.

Pf: $f(x_{t+1}) \leq f(x_t) - \frac{\mu}{2} \|\nabla f(x_t)\|_2^2$

sum over $t=1, \dots, T$

$$\sum_{t=1}^T f(x_t) \leq \sum_{t=0}^{T-1} f(x_t) - \frac{\mu}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

$$\Rightarrow f(x_T) \leq f(x_0) - \frac{\mu}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2$$

$$\Rightarrow \frac{\mu}{2} \sum_{t=0}^{T-1} \|\nabla f(x_t)\|_2^2 \leq f(x_0) - f(x_T)$$

$$\frac{\mu}{2} \cdot T \cdot \min_{0 \leq t \leq T-1} \|\nabla f(x_t)\|_2^2 \leq f(x_0) - \min_x f(x)$$

$$\min_{0 \leq t \leq T-1} \|\nabla f(x_t)\|_2 \leq \sqrt{\frac{f(x_0) - \min_x f(x)}{\mu(T-1)}} = o(\epsilon)$$

Gradient Descent for Quadratic Functions

Problem: $\min_x \frac{1}{2} x^T A x$ with $A \in \mathbb{R}^{d \times d}$ being positive-definite. $x = 0$

Theorem: Let λ_{\max} and λ_{\min} be the largest and the smallest eigenvalues of A . If we set $\eta \leq \frac{1}{\lambda_{\max}}$, we have

$$\begin{aligned} \|x_t\|_2 &\leq (1 - \eta \lambda_{\min})^t \|x_0\|_2 \\ \|x_{t+1}\|_2 &= \|x_t - \eta A x_t\|_2 \\ &= \|(I - \eta A) x_t\|_2 \\ &\leq \|I - \eta A\|_2 \|x_t\|_2 \\ &\leq (1 - \eta \lambda_{\min}) \|x_t\|_2 \\ &\leq (1 - \eta \lambda_{\min})^{t+1} \|x_0\|_2 \end{aligned}$$

To make $\|x_t\|_2 \leq \epsilon$

$$\eta \leq \frac{1}{\lambda_{\max}}$$

\Rightarrow need $\frac{\lambda_{\max}}{\lambda_{\min}} \log\left(\frac{1}{\epsilon}\right)$ iterations

$\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$ condition number

$$A = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & -\lambda_n \end{pmatrix}$$

Momentum: Heavy-Ball Method (Polyak '64)

$$V_0 = 0$$

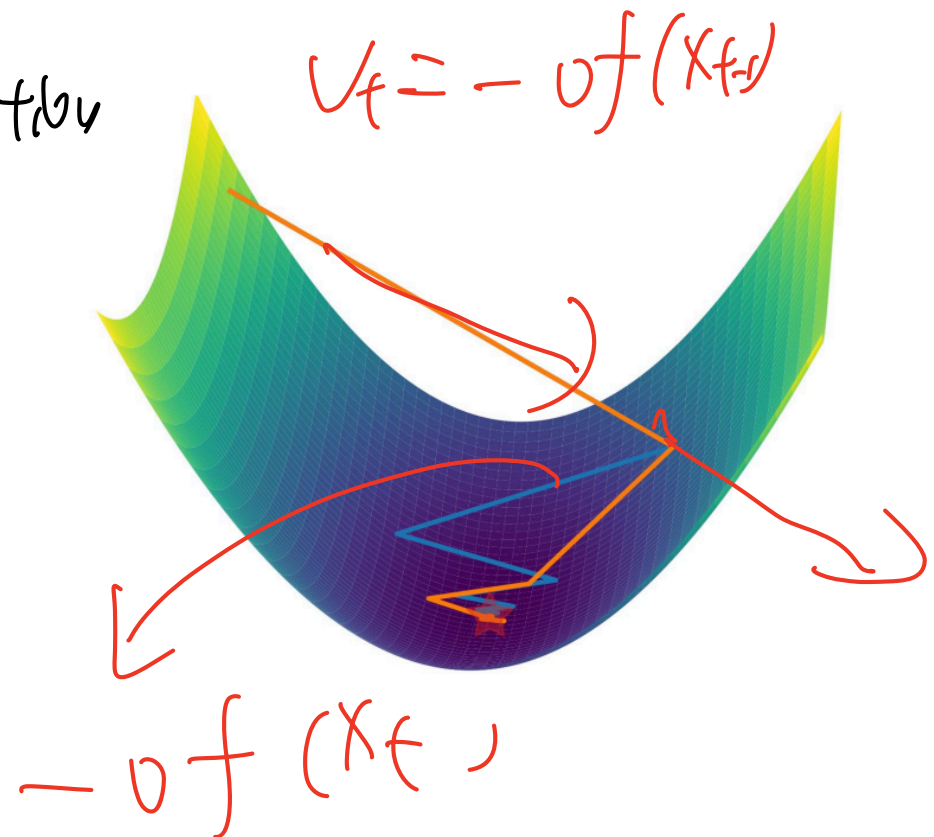
Problem: $\min_x f(x)$

Method: $v_{t+1} = -\nabla f(x_t) + \beta v_t$
 $x_{t+1} = x_t + \eta v_{t+1}$

for quadratic optimization

$$O\left(\sqrt{\kappa} \log\left(\frac{1}{\epsilon}\right)\right)$$

$$v.s. \quad \kappa \log\left(\frac{1}{\epsilon}\right)$$



Momentum: Nesterov Acceleration (Nesterov '89)

Problem: $\min_x f(x)$

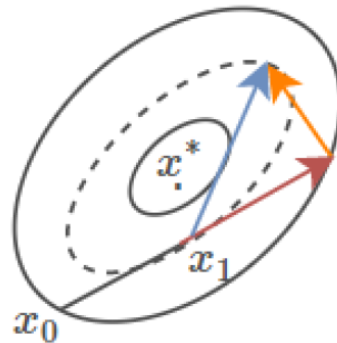
1. $\nabla f(x_t)$
2. $\nabla f(x_t + \beta v_t)$

THEORY

$\mathcal{O}(\log \frac{1}{\epsilon})$
for general
strongly convex
functions

Method: $v_{t+1} = -\nabla f(x_t + \beta v_t) + \beta v_t$
 $x_{t+1} = x_t + \eta v_{t+1}$

Polyak's Momentum



Nesterov Momentum

