

# Clarke Differential

---

**W**

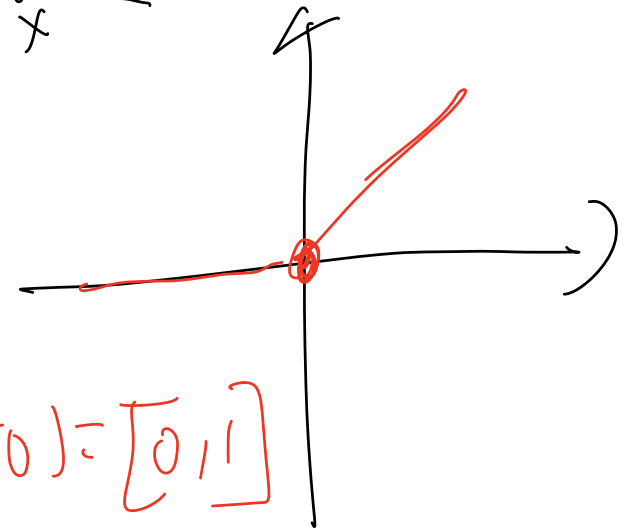
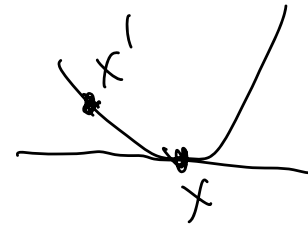
# Subdifferential and Subgradient

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

$\partial_s f(x) \triangleq \{s \in \mathbb{R}^d : \forall x' \in \mathbb{R}^d, f(x') \geq f(x) + \underbrace{s^\top (x' - x)}_{\text{linear function}}\}$ . The elements in the subdifferential set are subgradients.

$$x_{t+1} = x_t - \eta g_t$$

$$g_t \in \partial_s f(x)$$



$$\partial_s f(0) = [0, 1]$$

# Subdifferential and Subgradient

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

$\partial_s f(x) \triangleq \{s \in \mathbb{R}^d : \forall x' \in \mathbb{R}^d, f(x') \geq f(x) + s^\top(x' - x)\}$ . The elements in the subdifferential set are subgradients.

• If  $f$  is convex  $\Rightarrow \partial_s f(x)$  exists everywhere

• If  $f$  is convex and differentiable  $\partial_s f(x) = \{ \nabla f \}$  & smooth  $\Downarrow$

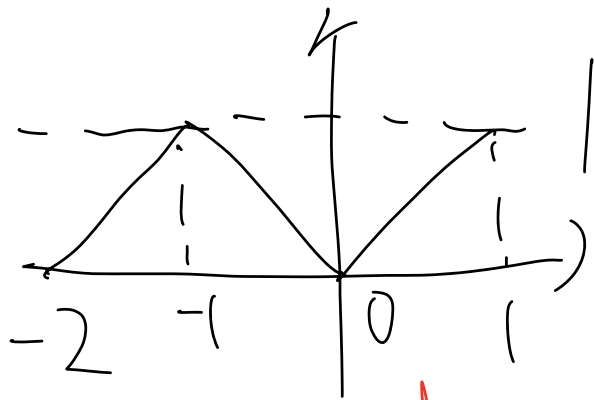
$\mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$  rate  $\mathcal{O}\left(\frac{1}{T}\right)$

# Subdifferential is not enough

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

$\partial_s f(x) \triangleq \{s \in \mathbb{R}^d : \forall x' \in \mathbb{R}^d, f(x') \geq f(x) + s^T(x' - x)\}$ . The elements in the subdifferential set are subgradients.

Problem:  $|\mathbb{N}|$  is not convex



$\Rightarrow$  well-defined

$x = -1$   
 we need  $s$  s.t.  $\forall x'$   
 $f(x') \geq f(-1) + s \cdot (x' - (-1))$

choose  $x' = -2$   
 $0 \geq 1 + s \cdot (-1) \Rightarrow s \geq 1$

choose  $x' = 1$   
 $\partial_s f(-1) = \emptyset$  s.t.  $1 \geq 1 + s \cdot (2) \Rightarrow s \leq 0$

# Clarke Differential

$$x_{t+1} = x_t - \eta g_t$$

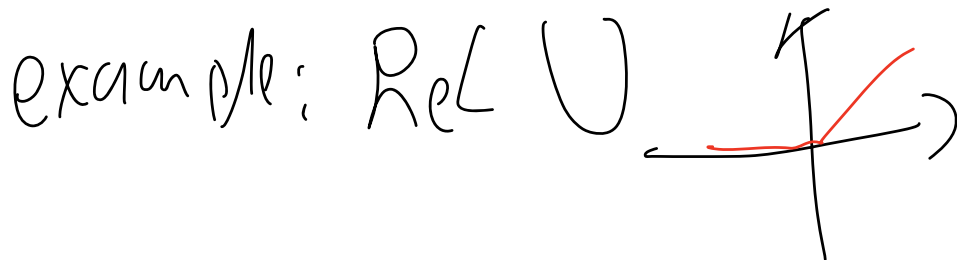
$$g_t \in \partial f(x)$$

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

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

The elements in the subdifferential set are subgradients.

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



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

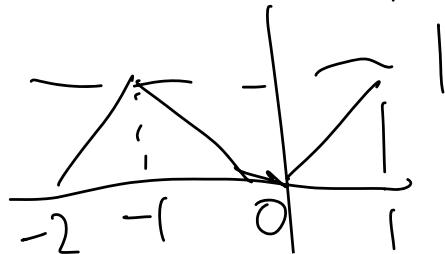
$$\partial f(x_i) = 0$$

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

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

$$\text{conv}(\{0, 1\}) = [0, 1]$$

example:



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

$$\{x_i\} : 0, -\frac{1}{2}, \dots \rightarrow -1, \partial f(x_i) = -1$$

$$\Rightarrow \text{conv}(\{1, -1\}) = [-1, 1]$$

# When does Clarke differential exists

$$x: \exists S, \forall x', |f(x) - f(x')| \leq L \|x - x'\|$$

**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  $S$ .

• If  $f$  is locally Lipschitz  $\rightarrow \partial f$  exists everywhere

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

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

# Positive Homogeneity

$$b(z) = \max\{z, 0\}$$

**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 ReLU:  $b(\alpha z) = \alpha \cdot b(z)$
- ② monomials of degree  $L$ :  $\prod_{i=1}^d x_i^{p_i}$ ,  $\sum_{i=1}^d p_i = L$
- $$\frac{d}{\prod_{i=1}^d (\alpha x_i)^{p_i}} = \alpha^{\sum p_i} \cdot \frac{d}{\prod_{i=1}^d x_i^{p_i}}$$
- $$= \alpha^L \frac{d}{\prod_{i=1}^d x_i^{p_i}}$$
- ③ Norm:  $\|\alpha x\| = \alpha \cdot \|x\|$

# Positive Homogeneity

④ Multi-layer ReLU

$$f(x, w_1, \dots, w_{H+1}) = w_{H+1} \sigma(w_H \dots \sigma(w_1 x) \dots)$$

— for one-layer

$$f(x, w_1, \dots, \alpha w_{H+1}, w_H, \dots, w_1, x) = \alpha \cdot w_{H+1} \sigma(w_H \dots \sigma(w_1 x) \dots)$$

— for all layers

$$f(x, \alpha w_1, \dots, \alpha w_{H+1}) = \alpha^{H+1} f(x, w_1, \dots, w_{H+1})$$

$\Rightarrow$   $(H+1)$ -homogeneous



# Positive Homogeneity

say  $W_h \in \mathbb{R}^{m \times m}$

Fact:  $\forall h = 1, \dots, H+1$   
 $\left\langle W_h, \frac{\partial f(x, w_1, \dots, w_{H+1})}{\partial w_h} \right\rangle = \underbrace{f(x, w_1, \dots, w_{H+1})}_{\text{independent of } h}$

Pf:  $A_h = \text{diag}(\sigma'(W_h x) \dots \sigma'(W_1 x) \dots) \in \mathbb{R}^{m \times m}$   
 $\sigma' = 0$  or  $1 \Rightarrow$  pattern whether activation is on or not

$$f(x, w_1, \dots, w_{H+1}) = W_{H+1} A_H W_H \dots A_1 W_1 x$$

$$\frac{\partial f}{\partial w_h} = (W_{H+1} A_H \dots W_{h+1} A_h)^T (A_{h-1} W_{h-1} \dots W_1 x)$$

$\Rightarrow$  verify

# Positive Homogeneity and Clark Differential

**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  $\langle s, x \rangle = Lf(x)$ .

$$\begin{aligned} \text{if } (100) \quad x &= W_n \\ s &: \frac{\partial f}{\partial W_n} \\ \langle s, x \rangle &= f(x) \end{aligned}$$

# Norm Preservation

$$f(x, w_1, w_2, w_3) = w_3 \sigma(w_2 \sigma(w_1 x))$$

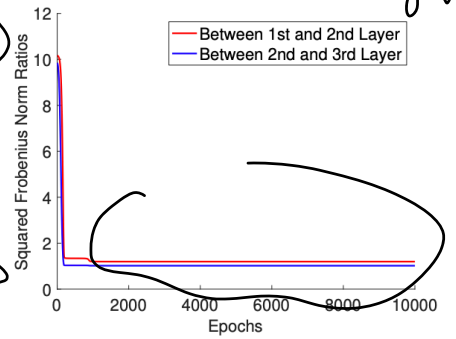
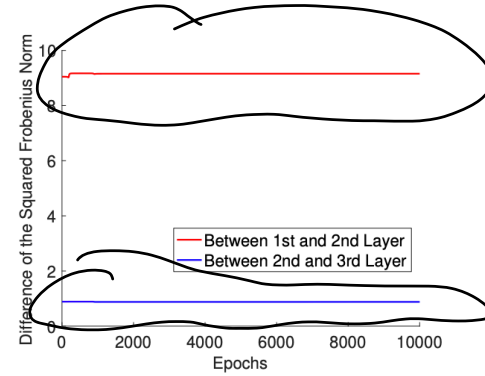
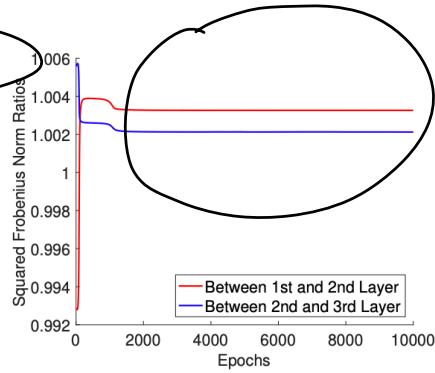
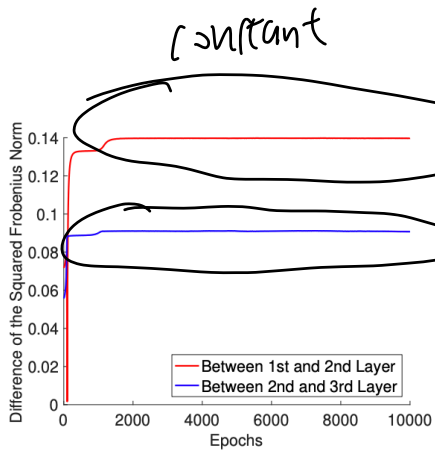
quadratic loss

$$w_3 \times 10$$

$$w_2 / 10$$

$$w_2(t+1) = w_2(t) - \eta \frac{\partial L}{\partial w_2}$$

$$w_3(t+1) = w_3(t) - \eta \frac{\partial L}{\partial w_3}$$



(a) Balanced initialization, squared norm differences.

(b) Balanced initialization, squared norm ratios.

(c) Unbalanced Initialization, squared norm differences.

(d) Unbalanced initialization, squared norm ratios.

(a) —  $\|w_1\|_F^2 - \|w_2\|_F^2$

(b) —  $\frac{\|w_1\|_F^2}{\|w_2\|_F^2}$

$\|w_1\|_F^2 - \|w_1(0)\|_F^2 = \|w_2\|_F^2 - \|w_2(0)\|_F^2$

if init small  $\Rightarrow$  balanced layers  
 $w_1(0), w_2(0)$  small  
 $\Rightarrow \|w_1(t)\|_F^2 \approx \|w_2(t)\|_F^2$

# Gradient flow and gradient inclusion

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

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

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

$$\frac{x_{t+1} - x_t}{\eta} = -\nabla f(x)$$

$$\text{let } \eta \rightarrow 0$$

# Norm preservation by gradient inclusion

**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  $t \geq 0$ ,

$$\frac{1}{2} \|W_{h_j}^c(t)\|_F^2 - \frac{1}{2} \|W_{h_j}^c(0)\|_F^2 = \frac{1}{2} \|W_{i_j}(t)\|_F^2 - \frac{1}{2} \|W_{i_j}(0)\|_F^2.$$

$\Rightarrow$   $\|W_i(t)\|_F^2 \approx \|W_j(t)\|_F^2$  if init is small

Proof:  $\int \frac{d \|W_i(t)\|_F^2}{dt} dt, \quad \frac{d \|W_i(t)\|_F^2}{dt} = 2 \langle W_i(t), \frac{dW_i(t)}{dt} \rangle = 2 \langle W_i(t), -\frac{dL}{dW_i} \rangle$

$\Rightarrow$  apply lemma  
 $\Rightarrow$  independent of  $i$ .