# **Clarke Differential**



#### Subdifferential and Subgradient

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

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

$$\mathbb{R}^{d}: \forall x' \in \mathbb{R}^{d}, f(x') \geq f(x) + s^{\mathsf{T}}(x' - x) \}. \mathsf{TI}$$
 subdifferential set are subgradients. 
$$\forall t = x + -y \neq t$$
 
$$\forall t \in \mathcal{I}(x)$$
 
$$\forall t \in$$

#### Subdifferential and Subgradient

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

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

$$O\left(\frac{1}{\sqrt{17}}\right)$$

#### Subdifferential is not enough

**Definition**: Given  $f: \mathbb{R}^d \to \mathbb{R}$ , 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.

ments in the subdifferential set are subgradients.

Roblem: No opt of ws (suvex)

$$x = -1$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

$$f(x) 7 | + S \cdot (x - (-1))$$

#### Clarke Differential

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

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

#### When does Clarke differential exists

bx' | f(x')-fk( ≤ L1 |(x'-x)(

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

### **Positive Homogeneity**

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

$$(\alpha x) = \alpha^{L} f(x) \text{ for any } \alpha \geq 0.$$

$$(1) \text{ ReLU: } 6(\lambda^{2}) = \lambda \cdot 6(\lambda^{2})$$

$$(2) \text{ monomials of degine } (\lambda^{2}) = \lambda^{2} \cdot 6(\lambda^{2})$$

$$(3) \text{ Novay } ||\lambda^{2}|| = \lambda^{2} \cdot ||\lambda^{2}||$$

$$(3) \text{ Novay } ||\lambda^{2}|| = \lambda^{2} \cdot ||\lambda^{2}||$$

### **Positive Homogeneity**

**Positive Homogeneity** 

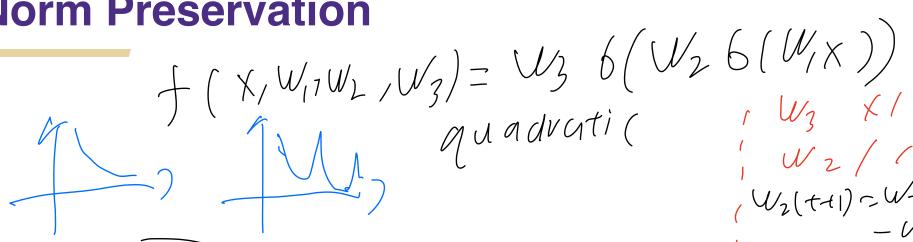
Fact 
$$\forall h = 1, ..., |H|/$$
 $\leq W_{h_1} \frac{\partial f(X, W_{1}, ..., W_{H+1})}{\partial W_{h_1}} = f(X, W_{1}, ..., W_{H+1})$ 

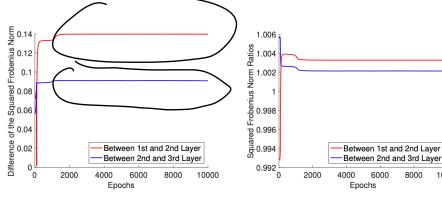
Of:  $A_{h_1} = Jiug(6'(W_{h_1}6(-...6(W_{1}X)-...))$ 
 $6' = 0$  or  $|$ 
 $f(X, W_{1}, ..., W_{H+1}) = W_{H+1}A_{H_1}W_{H_1} - ... A_{1}W_{1}X$ 
 $\frac{\partial f}{\partial W_{h_1}} = (W_{H+1}A_{H_1} ... W_{h_1}, A_{h_1})^T(A_{h_1} W_{h_1}... W_{h_1})$ 
 $= (W_{H+1}A_{H_1} ... W_{h_1}, A_{h_1})^T(A_{h_1} W_{h_1}... W_{h_1})$ 

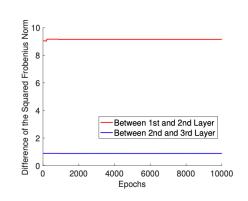
#### Positive Homogeneity and Clark Differential

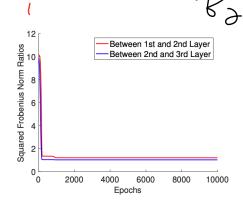
**Lemma:** Suppose  $f: \mathbb{R}^d \to \mathbb{R}$  is Locally Lipschitz and L -positively homogeneous. For any  $x \in \mathbb{R}^d$  and  $s \in \partial f(x)$ , we have  $\langle s, x \rangle = Lf(x)$ .

#### **Norm Preservation**









- Balanced initialization, squared norm differ-
- Balanced initialization, squared norm ratios.
- (c) Unbalanced Initialization, squared norm differences.
- (d) Unbalanced initialization, squared norm ratios.

 $F - ||W_2||_F^2 = ||W_0|||_F^2 - ||W_2(0)||_F^2$   $f = ||W_1(0)||_F - ||W_2(0)||_F^2$   $f = ||W_1(0)||_F - ||W_2(0)||_F^2$ 

#### Gradient flow and gradient inclusion

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

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

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

#### Norm preservation by gradient inclusion

**Theorem** (Du, Hu, Lee '18) Suppose  $\alpha > 0$ ,  $f(x; (W_{H+1}, ..., \alpha W_i, ..., W_1)) = \alpha f(x, (W_{H+1}, ..., W_1)), \text{ i.e.,}$ predictions are 1-homogeneous in each layer. Then for every pair of layers  $(I,I) \in [H+1] \times [H+1]$ , the gradient inclusion maintains: for all  $t \ge 0$ ,  $\frac{1}{2} \|W_h(t)\|_F^2 - \frac{1}{2} \|W_h(0)\|_F^2 = \frac{1}{2} \|W_h(t)\|_F^2 - \frac{1}{2} \|W_{h'}(0)\|_F^2.$  $=) \left( |W_{4}(t)||_{F}^{2} - |W_{4}(t)||_{F}^{2} + |W_{4}(0)||_{F}^{2} - |W_{4}(0)||_{F}^{2} \right)$ 0+: S d | | Wh (1) | | 2 d /

# Optimization Methods for Deep Learning



### Gradient descent for non-convex optimization

**Decsent Lemma:** Let  $f: \mathbb{R}^d \to \mathbb{R}$  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:

## **Converging to stationary points**

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

Pf:  $f(Xt) \le f(Xt) - \frac{1}{2} \| of(Xt)\|_2^2$ 

Sum over  $t = 0$ , ...,  $T - 1$ 
 $f(Xt) \le \frac{T-1}{2} | f(Xt) - \frac{1}{2} | f(Xt)|_2^2$ 

=)  $f(Xt) \le f(Xt) - \frac{1}{2} | f(Xt)|_2^2$ 

=)  $f(Xt) \le f(Xt) - \frac{1}{2} | f(Xt)|_2^2$ 

=)  $f(Xt) \le f(Xt) - \frac{1}{2} | f(Xt)|_2^2$ 

Theorem:  $f(Xt) = f(Xt) - \frac{1}{2} | f(Xt)|_2^2$ 

The