

# Lecture 20: Positive Homogeneity & low norm solutions

Reminders:

(Clarke differential)  $\partial f(\omega) = \left\{ \lim_{i \rightarrow \infty} Df(\omega_i) : \begin{array}{l} \omega_i \rightarrow \omega \\ Df(\omega_i) \text{ exists} \\ (\exists i, Df(\omega_i) \text{ exists}) \end{array} \right\}$

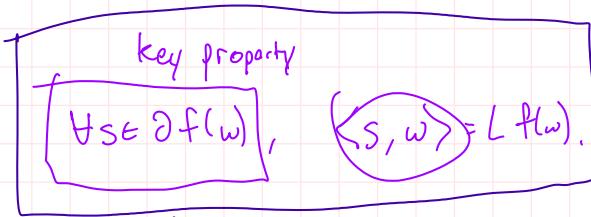
Good: locally Lipschitz  $\Rightarrow$  nonempty everywhere

Bad: technical issues (arc. regularity, nonuniqueness of flows)

does not match pytorch

Positive homogeneity:  $f$  is  $L$  homogeneous  $\forall c \geq 0 \quad f(cx) = c^L f(x)$ .

(e.g., abstraction of  $L$ -layer ReLU network without biases)



$$\sigma(cx) = \max\{0, cx\} = c\sigma(x)$$

$$\begin{aligned} F(x; cw) &= F(x; (cW_2, cb_2, cW_1)) \\ &= cW_2\sigma(cW_1x) + cb_2 \end{aligned}$$

$$= c^2W_2\sigma(W_1x) + \boxed{\begin{matrix} cb_2 \\ ? \end{matrix}} \quad \begin{array}{l} \text{not} \\ \text{pos-hom} \\ \text{if } b_2 \neq 0. \end{array}$$

Plan: ① prove & example

② "norm preservation"  $\|W_k(t) - W_k(0)\|$   
vs.  $\|W_k(t) - W_k(0)\|$

③ low norm solutions in optimization

## Positive homogeneity & gradients

$$\langle \partial f(w), w \rangle = \{ L f(w) \} \quad ???$$

ReLU network:

$$z_i := W_i x$$

$$z_{i+1} := W_{i+1} \sigma(z_i)$$

note

$$S := \text{diag}(\{1, 0\})$$

$$\sigma(z) = S z$$

$$F(x; w) = W_L \sigma(z_L) = W_L S_{L-1} z_{L-1} = W_L S_{L-1} W_{L-1} \sigma(z_{L-1})$$

$$= \dots = W_L S_{L-1} W_{L-1} S_{L-2} \dots S_1 w.$$

$$\left[ \frac{d}{d w_j} F(x; w) \right]_{w_j \in \mathbb{R}^{d_{L+1} d_L}} = (W_L S_{L-1} W_{L-1} \dots W_{j+1} S_j)^T \underbrace{(S_{j-1} W_{j-1} \dots S_1 w)}_{\text{dim } d_L}.$$

$$\langle \partial f(x; w), w \rangle = \sum_{j=1}^L \langle \partial_{w_j} f(x; w), w_j \rangle$$

memory cost  
 ① write out activation  
 ②  $\frac{d}{d x_i} x_1 x_2 \dots x_i x_{i+1} \dots x_n = x_1 \dots x_{i-1} e^{x_i} e^{x_{i+1}} \dots x_n$   
 ③ use transposes to match dimensions

$$\langle \partial_{w_i} F(x; w), w_i \rangle = \langle (W_L \dots S_i)^T (S_{i-1} \dots W_1)^T, w_i \rangle$$

$$= \text{tr}((W_L \dots S_i)^T (S_{i-1} \dots W_1)^T w_i)$$

$$= \text{tr}((S_{i-1} \dots W_1) (W_1 \dots S_i) w_i)$$

$$= \text{tr}(\nabla f(x; w)) = f(x; w).$$

$$\Rightarrow \langle \nabla F(x; w), w \rangle = L \cdot F(x; w)$$

Reminder to future self: instead:

① ReLU network is piecewise affine in the input  $x$ , piecewise polynomial in the parameters

② locally Lipschitz  $\Rightarrow$  nondifferentiability has Lebesgue measure zero

③ (i) calculate gradient exterior of each piece (ii) shifts towards boundaries.

Theorem (Euler's homogeneous function theorem)  
 If  $f$  is locally Lipschitz  $\Rightarrow \langle \nabla f(w), w \rangle = \{ L f(w) \}$

Proof. 2-phase proof: ① check when differentiable ② use limits & defn Clarke.

① When differentiable:

$$\begin{aligned} \langle \nabla f(w), w \rangle &= \lim_{h \rightarrow 0} \frac{f(w+hw) - f(w)}{h} \\ &= \lim_{h \rightarrow 0} \frac{f((1+h)w) - f(w)}{h} \\ &= \lim_{h \rightarrow 0} \frac{(1+h)^L f(w) - f(w)}{h} \\ &= \lim_{h \rightarrow 0} \frac{(1+Lh + \frac{1}{2}h^2 \dots) f(w) - f(w)}{h} \approx \lim_{h \rightarrow 0} (Lh f(w) + O(h^2 \dots)) \\ &= L f(w). \end{aligned}$$

$\min_w \sum_i \|w - x_i\|$ ,  
 Clarke book:  
 has drawing &  
 the table

$$\langle s, w \rangle = \sum_i x_i \langle s_i, w \rangle = \sum_i x_i \langle \lim_j \langle \partial f(w_{i,j}), w \rangle, w \rangle$$

$$= \sum_i x_i \lim_j \langle \partial f(w_{i,j}), w \rangle = \sum_i x_i \lim_j L f(w_{i,j})$$

$$= \sum_i x_i L f(w) = L f(w).$$

done (Clarke diff.).

if  $L \in \mathbb{R}_{\geq 0} \setminus \mathbb{Z} \Rightarrow$

$$\begin{aligned} &\lim_{h \rightarrow 0} \frac{f((1+h)w) - f(w)}{h} \\ &= \lim_{h \rightarrow 0} \frac{(1+h)^L f(w) - f(w)}{h} \\ &\stackrel{?}{=} \lim_{h \rightarrow 0} \frac{\exp(hL) f(w) - f(w)}{h} \end{aligned}$$

into the trash.

Theorem. Consider  $\nabla F(w_t) = -\sum_{i=1}^n l'(y_i; f(x_i; w_t))$   
 where  $F(x; w)$  is 1-homogeneous in each layer  
 $\sum_{i=1}^n F(x_i; (w_1, \dots, c w_i, w_{i+1}, \dots, w_n)) = c F(x_i; w_1, \dots, w_n)$

$$\text{Then } H_{i,k} \frac{1}{2} \|W_i(t)\|^2 - \frac{1}{2} \|W_i(0)\|^2 = \frac{1}{2} \|W_i(t)\|^2 - \frac{1}{2} \|W_i(0)\|^2.$$

Remark: for other notions of "distance to initialize" might be false.

$$\begin{aligned} \frac{1}{2} \|W_i(t) - W_i(0)\|^2 &= \frac{1}{2} \langle W_i(t), W_i(t) - W_i(0) \rangle \\ &= \frac{1}{2} \langle W_i(t), W_i(t) \rangle - \frac{1}{2} \langle W_i(t), W_i(0) \rangle \end{aligned}$$

$$= \int_0^t \frac{d}{ds} \frac{1}{2} \|W_i(s)\|^2 ds$$

$$= \int_0^t \langle \dot{W}_i(s), W_i(s) \rangle ds$$

$$= \int_0^t -\frac{1}{n} \sum_i l' \left( y_i; f(x_i; w_s) \right) \langle \dot{W}_i(s), F(x_i; w_s, W_i(s)) \rangle ds$$

$$= \int_0^t -\frac{1}{n} \sum_i l' \left( y_i; f(x_i; w_s) \right) F(x_i; w_s, W_i(s)) ds$$

no dependence on layer.

other note to

future self:

prove Clarke decomposition

over coordinates