

Cheatsheet for Fisher-Rao Metric, Geometry, and Complexity of Neural Networks

Sandro Braun, Leander Kurscheidt

July 29, 2018

Abstract

1 Geometry of Deep Rectified Networks

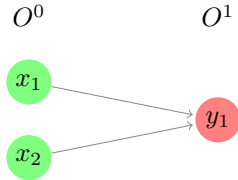
1.1 Lemma 2.1

Lemma 1.1 (Structure in Gradients).

$$\sum_{t=0}^L \sum_{i \in [k_t], j \in [k_{t+1}]} \frac{\partial O^{L+1}}{\partial W^{t_{ij}}} W_{ij}^t = (L+1)O^{L+1}(x) = \langle \nabla_{\theta} f_{\theta}(x), \theta \rangle \quad (1)$$

1.1.1 Example for Lemma 2.1

Take a simple feed forward neural network with only two inputs and one output. Remember that by definition $O^0 = x$. Illustrated:



The simple example can then be written with:

$$O^0 = \begin{bmatrix} O_1^0 \\ O_2^0 \end{bmatrix}, W^0 = \begin{bmatrix} W_1^0 \\ W_2^0 \end{bmatrix} \quad (2)$$

$$O^1 = \sigma(O^{0T} W^0) = \sigma(\underbrace{O_1^0 W_1^0 + O_2^0 W_2^0}_z). \quad (3)$$

Substituting the term in brackets will simplify notation.

Now calculate the partial derivatives of (2) with respect to W_i^0 .

$$\begin{aligned} \frac{\partial O^1}{\partial W_1^0} &= \sigma'(z) O_1^0 \\ \frac{\partial O^1}{\partial W_2^0} &= \sigma'(z) O_2^0 \end{aligned}$$

Summing up the $\frac{\partial O^1}{\partial W_j^0} W_j^0$ reveals,

$$\begin{aligned} \sum_{j=1}^2 \frac{\partial O^1}{\partial W_j^0} W_j^0 &= \sigma'(z) O_1^0 W_1^0 + \sigma'(z) O_2^0 W_2^0 = \sigma'(z) (\underbrace{O_1^0 W_1^0 + O_2^0 W_2^0}_z) \\ &= \sigma'(z) z \end{aligned}$$

Note that this is equivalent to calculating $\langle \nabla_{W^0} O^1, W^0 \rangle$:

$$\nabla_{W^0} O^1 = \begin{bmatrix} \frac{\partial O^1}{\partial W_1^0} \\ \frac{\partial O^1}{\partial W_2^0} \end{bmatrix}$$

$$\langle \nabla_{W^0} O^1, W^0 \rangle = \sum_{j=1}^2 \frac{\partial O^1}{\partial W_j^0} W_j^0 = \dots = \sigma'(z)z$$

Using the relation $\sigma(z) = \sigma'(z)(z)$ reveals (2). Then reading left to right reveals (1) for $L = 0$, $\theta = W^0$ and $f_\theta(x) = O^1(x)$ which completes the example.

$$\sum_{t=0}^{L=0} \sum_{j=1}^2 \frac{\partial O^1}{\partial W_j^t} W_j^t = \langle \nabla_{W^0} O^1, W^0 \rangle = \sigma'(z)z = \sigma(z) = O^1$$

1.2 Corollary 2.1

1.2.1 Notes

1. *Proof.* We want to show $\frac{\partial l(f,Y)}{\partial f} = -y \Leftrightarrow yf < 1$. So,

$$\begin{aligned} & 1 - y_i f_i > 0 \\ \Leftrightarrow & l = 1 - y_i f_i \\ \Leftrightarrow & \frac{\partial l}{\partial f} = -y_i \end{aligned}$$

□

2. *Proof.* We want to show $\frac{\partial l(f,Y)}{\partial f} = 0 \Leftrightarrow yf > 1$. So,

$$\begin{aligned} & 1 - y_i f_i < 0 \\ \Leftrightarrow & l = 0 \\ \Leftrightarrow & \frac{\partial l}{\partial f} = 0 \end{aligned}$$

□

References

$t=1, s=1 \quad 0 \leq t \leq s \leq L$

$O^{t+2}(x) = \sigma(\underbrace{O_1^{t+1} w_1 + O_2^{t+1} w_2}_z)$

$\frac{\partial O^2}{\partial w_1} = \sigma'(z) \cdot O_1^1 \cdot w_1$

$\frac{\partial O^1}{\partial w_2} = \sigma'(z) \cdot O_2^1 \cdot w_2$

$= \sigma'(z) \left(\underbrace{O_1^1 \cdot w_1 + O_2^1 \cdot w_2}_z \right)$

$= \sigma(z)$

$f_\theta(x) = \sigma(x^T \cdot \theta)$

$\frac{\partial f}{\partial \theta} = \sigma'(x \cdot \theta) \cdot x$

$\left\langle \frac{\partial L}{\partial \theta^1}, \theta \right\rangle = \sigma'(x \cdot \theta) \cdot x \cdot \theta$

$\sigma'(x \cdot \theta)$

$= f_\theta(x)$

$0 \quad 1 \quad 2 \quad \dots \quad L$

$\uparrow \log t$

$N^{t,m} \rightarrow O^t(x)$

$O^t(y) = \sigma_t(N^t(y))$

non-linear

$X \in \mathbb{R}^{2 \times 2} \rightarrow Y \in \mathbb{R}^{1 \times 1}$

$w^0 \in \mathbb{R}^2, w_1^2, w_2^2, w_3^2 = 10$

$2.2: \sigma(x) = \sigma'(x) \cdot x$

$O^0(x) = x \in \mathbb{R}^p$

$O^{L+1}(x) = f_\theta(x) \in \mathbb{R}^k$

$\sigma \left(\underbrace{\sigma(x^T w^0)}_{n^1} \right)$

n^2

$$\ell = \max_{\gamma} \{0, \underbrace{1 - \gamma_i f_i}_{\gamma_i}\}$$

1) $\frac{\partial \ell(f, \gamma)}{\partial f} = -\gamma \Leftrightarrow \gamma f < 1$

$\rightarrow 1 - \gamma_i f_i > 0 \Rightarrow \ell = 1 - \gamma_i f_i$

$\rightarrow \frac{\partial \ell}{\partial f} = -\gamma_i \quad \text{q.e.d.}$

2) $\gamma f \geq 1 \Rightarrow 1 - \gamma_i f_i \leq 0 \Rightarrow \ell = 0$

$\frac{\partial \ell}{\partial f} = 0 \quad \text{q.e.d.}$

1. $\partial_{\theta} \hat{\ell}(\theta) = \ell$

2. $\gamma_i f_{\theta}(x_i) \geq 0 \quad \forall i$

\Updownarrow

$\gamma_i f_{\theta}(x_i) \geq 1 \quad \forall i$

P.5 | Proof of Corollary 2.1

$$\hat{\ell}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(x_i), \gamma_i)$$

$$\langle \nabla_{\theta} \hat{\ell}(\theta), \theta \rangle = \frac{1}{N} \sum_{i=1}^N \langle \nabla_{\theta} \ell(\underbrace{f_{\theta}(x_i)}_{\theta^{L+1}(x_i)}, \gamma_i), \theta \rangle$$

$$= \frac{1}{N} \sum_{i=1}^N \ell'(f_{\theta}(x_i), \gamma_i) \langle \nabla_{\theta} f_{\theta}(x_i), \theta \rangle$$

$$= \frac{1}{N} \sum_{i=1}^N \ell'(f_{\theta}(x_i), \gamma_i) \cdot \underbrace{(L+1) \cdot f_{\theta}(x_i)}_{(2.6) + \nabla_{\theta}(\text{all } \theta \nabla)} \quad (2.6) + \nabla_{\theta}(\text{all } \theta \nabla)$$

$$= (L+1) \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell(f_{\theta}(x_i), \gamma_i)}{f_{\theta}(x_i)} f_{\theta}(x_i)$$

$$= (L+1) \hat{\mathbb{E}} \left[\frac{\partial \ell(f_{\theta}(x_i), \gamma_i)}{f_{\theta}(x_i)} f_{\theta}(x_i) \right]$$

1. $\partial_{\theta} \hat{\ell}(\theta) = \ell$

2. $\gamma_i f_{\theta}(x_i) \geq 0 \quad \forall i$

\Updownarrow

$\gamma_i f_{\theta}(x_i) \geq 1 \quad \forall i$

P.5 | Proof of Corollary 2.1

$$\hat{\ell}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(x_i), y_i)$$

$$\langle \nabla_{\theta} \hat{\ell}(\theta), \theta \rangle = \frac{1}{N} \sum_{i=1}^N \langle \nabla_{\theta} \ell(f_{\theta}(x_i), y_i), \theta \rangle$$

$$= \frac{1}{N} \sum_{i=1}^N \ell'(f_{\theta}(x_i), y_i) \langle \nabla_{\theta} f_{\theta}(x_i), \theta \rangle$$

$$= \frac{1}{N} \sum_{i=1}^N \ell'(f_{\theta}(x_i), y_i) \cdot (L+1) f_{\theta}(x_i) \quad (2.6) + \nabla_{\theta} (\text{all } \theta \nabla)$$

$$= (L+1) \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell(f_{\theta}(x_i), y_i)}{f_{\theta}(x_i)} f_{\theta}(x_i)$$

$$= (L+1) \hat{\mathbb{E}} \left[\frac{\partial \ell(f_{\theta}(x), y)}{f_{\theta}(x)} f_{\theta}(x) \right]$$

1. $\nabla_{\theta} \hat{\ell}(\theta) = \nabla$
2. $y_i, f_{\theta}(x_i) \geq 0 \forall i$
 \Uparrow
 $y_i, f_{\theta}(x_i) \geq 1 \forall i$

(3.1) $\|\theta\|_{\Gamma}^2 := \langle \theta, \mathbf{I}(\theta) \theta \rangle$, with $\mathbf{I}(\theta) = \mathbb{E} \left[\nabla_{\theta} \ell(f_{\theta}(x), y) \otimes \nabla_{\theta} \ell(f_{\theta}(x), y) \right]$

with $\mathbf{I}(\theta) = \mathbb{E} \left[\nabla_{\theta} \ell(f_{\theta}(x), y) \otimes \nabla_{\theta} \ell(f_{\theta}(x), y) \right]$

$\langle \theta, \nabla_{\theta} \ell(f_{\theta}(x), y) \rangle^2 = \langle \theta, \nabla_{\theta} \ell(f_{\theta}(x), y) \rangle \cdot \langle \theta, \nabla_{\theta} \ell(f_{\theta}(x), y) \rangle$

$= \left[\theta^T \nabla_{\theta} \ell(f_{\theta}(x), y) \right] \cdot \left[\nabla_{\theta} \ell(f_{\theta}(x), y)^T \theta \right]$

$= \theta^T \left(\underbrace{\nabla_{\theta} \ell(f_{\theta}(x), y) \nabla_{\theta} \ell(f_{\theta}(x), y)}_{\mathbf{I}} \right) \theta$

$\langle x, y \rangle = x^T y$

$$(3.7) \quad \|\theta\|_{fr}^2 := \langle \theta, I(\theta) \theta \rangle,$$

with $I(\theta) = E \left[\nabla_{\theta} l(f_{\theta}(x), y) \otimes \nabla_{\theta} l(f_{\theta}(x), y) \right]$

$T\theta_1 \theta_2 \left[\begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix} \right]$
 $m_1^2 \sigma_1^2 + m_2^2 \sigma_2^2$
 $+(m_{11}m_{21} + m_{12}m_{22})\sigma_1\sigma_2$

$$\begin{aligned} & \langle \theta, \nabla_{\theta} l(f_{\theta}(x), y) \rangle \\ &= \left\langle \left[\nabla_{\theta} l(f_{\theta}(x), y) \right]^T, \theta \right\rangle \end{aligned}$$

$$\stackrel{\text{chain}}{=} \left\langle \frac{\partial l(f_{\theta}(x), y)}{\partial f_{\theta}(x)} \nabla_{\theta} f_{\theta}(x), \theta \right\rangle$$

$$\begin{aligned} &= \left\langle \nabla_{\theta} f_{\theta}(x) \frac{\partial l(f_{\theta}(x), y)}{\partial f_{\theta}(x)}, \theta \right\rangle \\ &= \left\langle \frac{\partial l(f_{\theta}(x), y)}{\partial f_{\theta}(x)}, \nabla_{\theta} f_{\theta}(x)^T \theta \right\rangle \end{aligned}$$

$$\begin{aligned} \langle X, y \rangle &= x^T y \\ \langle \alpha x, y \rangle &= \alpha \langle x, y \rangle \\ \langle \alpha x, y \rangle &= x^T (\alpha y) \end{aligned}$$

$$\text{Es war: } \langle \nabla_{\theta} f_{\theta}(x), \theta \rangle = (L+1) \theta^{L+1}(x) = (L+1) f_{\theta}(x)$$

$$E \left[\left\langle \frac{\partial l(f_{\theta}(x), y)}{\partial f_{\theta}(x)}, (L+1) f_{\theta}(x) \right\rangle^2 \right]$$

$$= (L+1)^2 E \left[\left\langle \frac{\partial l(f_{\theta}(x), y)}{\partial f_{\theta}(x)}, f_{\theta}(x) \right\rangle^2 \right] = \|\theta\|_{fr}^2$$

$$l = \frac{1}{2} \int_0^1 \frac{\partial g}{\partial t} = 1 \quad \Rightarrow \quad E \left[\left\langle f_{\theta}(x) - y, f_{\theta}(x) \right\rangle^2 \right]$$

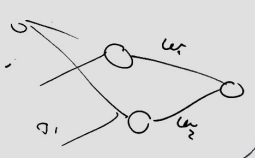
$$\nabla_{\theta} \theta^{L+1} = \begin{bmatrix} \frac{\partial \theta^{L+1}}{\partial \theta_1} \\ \vdots \\ \frac{\partial \theta^{L+1}}{\partial \theta_n} \end{bmatrix} \quad G = \begin{bmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{bmatrix}$$

$$f_{\theta}^2(x) - y f_{\theta}(x)$$

$$(\nabla \theta)^T G \nabla \theta$$

$$\nabla_{\Theta} O^{L+1}(x)^T \Theta$$

$$= \sum_{i=0}^L \frac{\partial O^{L+1}(x)}{\partial W_i} W_i$$

$$= \sum_{t=0}^L \sum_{i \in [k_t], j \in [k_{t+1}]} \frac{\partial O^{L+1}(x)}{W_{ij}^+} W_{ij}^+$$


$$\|G\|_{tr}^2 = \langle G, I(G)G \rangle$$

$$= (L+1)^2 E \left(\left\langle \frac{\partial Q(x,y)}{\partial t_G(x)}, t_G(x) \right\rangle^2 \right)$$

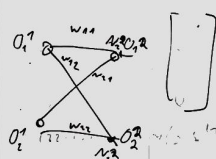
Norm von G ist unabhängig von G

$$(t-y)^2$$

$$t-y$$

Proof of 3.4

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

$$z = \begin{bmatrix} N_1 \\ N_2 \end{bmatrix}$$


$$\begin{bmatrix} \sigma'(N_1) N_1 \\ \sigma'(N_2) N_2 \end{bmatrix}$$

$$Zurgen: \underline{O}^2 = N^2 \text{diag}(\sigma'(N^2))$$

$$\underline{O}^1 = \begin{bmatrix} \sigma^1 \\ \sigma^2 \end{bmatrix} = \begin{bmatrix} \sigma(N_1) \\ \sigma(N_2) \end{bmatrix} = \begin{bmatrix} \sigma'(N_1) N_1 \\ \sigma'(N_2) N_2 \end{bmatrix} = \begin{bmatrix} N_1^2 & N_2^2 \end{bmatrix}$$

$$\begin{bmatrix} \sigma'(N_1) & 0 \\ 0 & \sigma'(N_2) \end{bmatrix}$$

$$\sigma'(N^2) = \begin{bmatrix} \sigma'(N_1) \\ \sigma'(N_2) \end{bmatrix}$$

Q.E.D.

$$\sigma(\sigma(x^T w_0) w_0)$$

$$(l+1)^2 E \left[\left(\frac{\partial f}{\partial t_0} \right)^2 \right]$$

$$f_1 = f_0(x) - \frac{\partial f}{\partial t_0} \bigg|_{t_0} t_0$$

$$f_0 = x^T w_0 D^T(x) \dots$$

$$(x^T v)^2 = (x^T v)(x^T v)$$

$$= (v^T x)(x^T v)$$

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Richtungsableitung 3.6
 $f(x, y | \theta) = p_\theta(x, y)$
 $\langle \nabla_\theta p_\theta(x), \alpha \rangle = \frac{d p_\theta(x)}{d t} = \bar{\alpha}$
 $\bar{\alpha} = \langle \nabla_\theta p_\theta(x), \alpha \rangle = (\nabla_\theta p_\theta(x))^T \alpha$

$$\langle \bar{\alpha}, \bar{\beta} \rangle = E_{(x, y) \sim p_\theta} \left[\frac{\partial}{\partial \theta} \frac{\beta}{p_\theta} \right] = E_{(x, y) \sim p_\theta} \left[\frac{(\nabla_\theta p_\theta(x))^T \alpha}{p_\theta(x)} \frac{(\nabla_\theta p_\theta(x))^T \beta}{p_\theta(x)} \right]$$

$$= E_{(x, y) \sim p_\theta} \left[\alpha^T (\nabla_\theta \log p_\theta(x)) \nabla_\theta \log p_\theta(x)^T \beta \right]$$

$$\langle \bar{\alpha}, I(\theta) \bar{\beta} \rangle = \alpha^T E_{(x, y) \sim p_\theta} [\nabla_\theta \log p_\theta(x) \nabla_\theta \log p_\theta(x)^T] \beta$$

$$\mathbb{R}^p \times \mathbb{R}^4$$

$$\nabla f(x, y)$$

$$\int x f(x) dx$$

4

-

1

11

U.S.

$$\int_{-\infty}^{\infty} x p_0(x) dx = E[X]$$

$$\int \frac{\bar{\alpha}}{p_{\alpha}} \cdot \frac{\bar{\beta}}{p_{\beta}} p_{\theta}(x) d(x, y) d(y)$$

EE

$$\frac{f'}{f} = (\log f)'$$

$$f = P_\theta(x) \Rightarrow \frac{\nabla_x P_\theta(x)}{P_\theta(x)} = \nabla_{\theta} \log P_\theta(x)$$