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

Sandro Braun, Leander Kurscheidt

July 25, 2018

**Abstract**

# 1 Geometry of Deep Rectified Networks

## 1.1 Lemma 2.1

**Lemma 1.1** (Structure in Gradients)**.**

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

### 1.1.1 Example for Lemma 2.2

$$\frac{\partial O^2}{\partial W_1} = \sigma'(z)O^1 W_1$$

$$\frac{\partial O^2}{\partial W_2} = \sigma'(z)O^1 W_2$$

therefore:

$$\sum_{t=1}^{2} \frac{\partial O^2}{\partial W_t} = \sigma'(z)(\underbrace{O^1 W_1 + O^1 W_2}_{z})$$

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



## 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,

$$1 - y_i f_i > 0$$
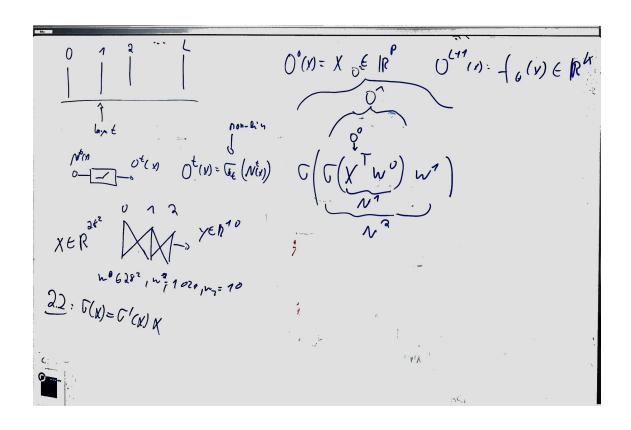$$\Leftrightarrow \quad l = 1 - y_i f_i$$
$$\Leftrightarrow \quad \frac{\partial l}{\partial f} = -y_i$$

$\square$

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

$$1 - y_i f_i < 0$$
$$\Leftrightarrow \quad l = 0$$
$$\Leftrightarrow \quad \frac{\partial l}{\partial f} = 0$$

$\square$

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

$O^{s+2}_{(x)} = \sigma\left( \overbrace{O_1^1 w_1 + O_2^1 w_2} \right)$

$z$

$k_i^{20}$

$O \quad w_1$
$\quad w_2 \quad O$
$O$
$O^s$

$O^{s+1}$

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

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

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

$= \sigma(z)$

$t_\theta(x) = \sigma\left( x^T \cdot \theta \right)$

$\dfrac{\partial t}{\partial \theta} = \sigma'(x \cdot \theta) \cdot x \quad \sigma(x \cdot \theta)$

$\left\langle \dfrac{\partial t}{\partial \theta}, \theta \right\rangle = \sigma'(x \cdot \theta) x \theta$

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

$= t_\theta(x)$

---

$0 \quad 1 \quad 2 \quad \cdots \quad L$

layer $t$

$N^{t}(x)$
$O \to \boxed{\phantom{x}} \to O^t(x)$

$O^t(x) = \sigma_t\left( N^t(x) \right)$

non-lin

$O^0(x) = x \quad O \in \mathbb{R}^p \quad O^{L+1}(x) = t_\theta(v) \in \mathbb{R}^k$

$O^1$

$O^0$

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

$N^2$

$X \in \mathbb{R}^{28^2} \quad \bowtie\bowtie \to Y \in \mathbb{R}^{10}$

$0 \quad 1 \quad 2$

$w^0 \; 28^2, \; w^1 \; 1024, \; w_2 = 10$

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

$$\dot{\ell} = \max\{0, \underbrace{1 - y_i f_i}_{\gamma_i}\}$$

1) $\dfrac{\partial \ell(f, Y)}{\partial f} = -y \quad \Leftrightarrow \quad yf < 1$

$\longrightarrow 1 - y_i f_i > 0 \Rightarrow \ell = 1 - y_i f_i$

$\longrightarrow \dfrac{\partial \ell}{\partial f} = -y_i \quad \text{q.e.d.}$

2) $yf \geq 1 \Rightarrow 1 - y_i f_i \leq 0 \Rightarrow \ell = 0$

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

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

$\Updownarrow$

$y_i f_\theta(x_i) \geq 1 \quad \forall i$

---

P.5 | Proof of Corollary 2.1

$\hat{L}(\theta) = \dfrac{1}{N} \sum_{i=1}^{N} \ell(f_\theta(x_i), Y_i)$

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

$= \dfrac{1}{n} \sum_{i=1}^{N} \ell'(f_\theta(x_i), Y_i) \langle \underbrace{\nabla_\theta f_\theta(x_i)}_{(L+1) \cdot f_\theta(x_i) \quad (2.6)} , \theta \rangle \; + \nabla_\theta (\text{all } \theta \; ?)$

$= \dfrac{1}{n} \sum_{i=1}^{N} \ell'(f_\theta(x_i), Y_i) \cdot (L+1) f_\theta(x_i)$

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

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

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

$\Updownarrow$

$y_i f_\theta(x_i) \geq 1 \quad \forall i$

| Proof of Corollar 2.1

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

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

$$= \frac{1}{n} \sum_{i=1}^{N} \ell'\left(f_\theta(x_i), y_i\right) \langle \underbrace{\nabla_\theta f_\theta(x_i)}_{(L+1) \cdot f_\theta(x_i) \quad (2.6) \; + \nabla_\theta \; (\text{all } \theta \, ?)}, \theta \rangle$$

$$= \frac{1}{n} \sum_{i=1}^{N} \ell'\left(f_\theta(x_i), y_i\right) \cdot (L+1) f_\theta(x_i)$$

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

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

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

$$\Updownarrow$$

$y_i \, f_\theta(x_i) \geq 1 \quad \forall i$

---

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

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

$\top \theta_1 \; \theta_2 \; \rbrack \begin{bmatrix} m_{11} & m_{12} \\ m_{21} & m_{21} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$

$m_{11}\theta_1^2 + \theta_2^2 m_{22}$
$+ (m_{12} + m_{21})\theta_1 \theta_2$

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

$\langle x, y \rangle := x^\top y$

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

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

$$(3.1) \qquad \|\theta\|_{fr}^2 := \langle \theta, I(\theta)\theta \rangle,$$

$$\text{with} \quad I(\theta) = E\left[\nabla_\theta l\big(f_\theta(X), Y\big) \otimes \underbrace{\hspace{4cm}} \right]$$

$$T\theta_1 \ \theta_2 \ \Big]\begin{bmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{bmatrix}$$

$$m_{11}\theta_1^2 + \sigma_2^2 m_{22}$$
$$+ (m_{12}+m_{21})\theta_1$$

$$\langle \theta, \nabla_\theta l\big(f_\theta(X), Y\big) \rangle$$
$$= \langle \nabla_\theta l\big(f_\theta(X), Y\big), \theta \rangle$$

$$\langle x, y \rangle = x^T y$$
$$\langle \alpha x, y \rangle = \alpha(x^T y)$$
$$= x^T(\alpha y)$$

$\overset{\text{chain}}{=} \langle \dfrac{\partial l(f_\theta(X), Y)}{\partial f_\theta(X)} \nabla_\theta f_\theta(X), \theta \rangle$

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

$$= \langle \nabla_\theta f_\theta(X) \, \partial \cdots, \theta \rangle$$

$$= \langle \dfrac{\partial l(f_\theta(X), Y)}{\partial f_\theta(X)}, \nabla_\theta f_\theta(X)^T \theta \rangle$$

s

---

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

$$E\left[ \langle \dfrac{\partial l(f_\theta(X,Y))}{\partial f_\theta(Y)}, (L+1) f_\theta(X) \rangle^2 \right]$$

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

$$l = \left(\dfrac{\ }{2}\right)^2 \to \dfrac{\partial l}{\partial f} = | \ | = 0 \quad E\left[ \langle f_\theta(X) - Y, f_\theta(X) \rangle^2 \right]$$

$$\dfrac{(f_\theta(X) - y)^2}{2} \qquad G = \begin{bmatrix} w_1 \\ \vdots \end{bmatrix} \qquad f_\theta^2(X) - Y f_\theta(X)$$

$$\nabla_\theta \theta^{L+1} = \begin{bmatrix} \partial \theta_{L+1} \\ \vdots \\ \partial u_n \end{bmatrix}$$

$$\nabla_\theta^T \theta = \qquad (\nabla \theta)^T G \theta$$

s

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

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

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



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

$$\vdots$$

$$= (L+1)^2 E\left( \langle \frac{\partial \ell(f(x;y))}{\partial f_\Theta(N)}, f_\Theta(x) \rangle^2 \right)$$

Norm
von
$\Theta$

ist unabhängig von $G$

$$\left(f - y\right)^2$$

$$f - y$$

---

Proof of 3.4

$$O(z) = \sigma'(z) z \qquad z = \begin{bmatrix} N_1 \\ N_2 \end{bmatrix}$$



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

zusammen: $\dot{O}^2 = \underline{N}^2 \, diag\left(\sigma'(\underline{N}^2)\right)$

$$\underline{O}^2 = \begin{bmatrix} O_1^2 \\ O_2^2 \end{bmatrix} = \begin{bmatrix} G(N_1^2) \\ \sigma(N_2^2) \end{bmatrix} = \begin{bmatrix} \sigma'(N_1^2) N_1^2 \\ \sigma'(N_2^2) N_2^2 \end{bmatrix} = \begin{bmatrix} N_1^2 & N_2^2 \end{bmatrix} \begin{bmatrix} \sigma'(N_1^2) & 0 \\ 0 & \sigma'(N_2^2) \end{bmatrix}$$

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

q.e.d

$y_i^A \in \mathbb{R}^P$

$$\underline{X}^T = \begin{bmatrix} P \\ \hline x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} N$$

Netzwerk hat skalaren output

$$\underline{X}^T \underline{V} \qquad \begin{bmatrix} x_1^T V \\ x_2^T V \\ \vdots \\ x_N^T V \end{bmatrix} \overset{1 \times 1}{=} \qquad N$$

$\Rightarrow V^T X = \boxed{N} \Rightarrow V^T X X^T V = \boxed{N} \; N = \boxed{1 \times 1}$

---

$(\log f)' = \dfrac{1}{f} \cdot f' = \dfrac{f'}{f} \qquad (-1)$

$f = P_{\theta + t\alpha} \Big|_{t=0} \qquad (\;)^t = \dfrac{d}{dt} \quad -1 \qquad es\ gilt: \; \dfrac{d P_{\theta + t\alpha}}{dt}\Big|_{t=0} = \overline{\alpha}$

$(1)$

$\Rightarrow \dfrac{\frac{d P_{\theta + t\alpha}}{dt}}{P_{\theta + t\alpha}}\Big|_{t=0} = \dfrac{\overline{\alpha}}{P_\theta}\Big| = \dfrac{f'}{f} = \dfrac{d}{dt}\left(\log P_{\theta + t\alpha}\right)\Big|_{t=0}$

$\Rightarrow \langle \overline{\alpha}, \overline{\beta} \rangle_{P_\theta} = \displaystyle\int_M \dfrac{\overline{\alpha}}{P_\theta} \dfrac{\overline{\beta}}{P_\theta} P_\theta = \displaystyle\int_{M = X^T} \dfrac{d}{dt} \log P_{\theta + t\alpha}\Big|_{t=0} \dfrac{d}{dt} \log P_{\theta + t\beta}\Big|_{t=0}$

Erwartungswert $= E\Big\{ \log P_{\theta + t\alpha} \; \log P_{\theta + t\beta} \Big\}$

Ableitung

# References