

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

Sandro Braun, Leander Kurscheidt

July 19, 2018

**Abstract**

## 1 Geometry of Deep Rectified Networks

Lemma 2.1 (Structure in Gradient)

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

## References



$$\ell = \max_{\gamma} \{0, 1 - \gamma_i t_i\}$$

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

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

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

2)  $\gamma t \geq 1 \Rightarrow 1 - \gamma_i t_i \leq 0 \Rightarrow \ell = 0$

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

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

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

$\Updownarrow$

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

P.5 | Proof of Corollary 2.1

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

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

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

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

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

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

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

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

$\Updownarrow$

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

P.5 | Proof of Corollary 2.1

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

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

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

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

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

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

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

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

$\begin{bmatrix} \theta_1 & \theta_2 \end{bmatrix} \begin{bmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$   
 $m_{11}\theta_1^2 + 2m_{12}\theta_1\theta_2 + m_{22}\theta_2^2$   
 $\langle x, y \rangle = x^T y$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

$$\nabla_{\theta} \phi^{L+1} = \begin{bmatrix} \frac{\partial \phi^{L+1}}{\partial \theta_1} \\ \vdots \\ \frac{\partial \phi^{L+1}}{\partial \theta_n} \end{bmatrix} \quad G = \begin{bmatrix} m_{11} & \dots & m_{1n} \\ \vdots & \ddots & \vdots \\ m_{n1} & \dots & m_{nn} \end{bmatrix}$$

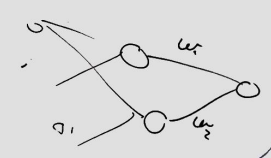
$$f_{\theta}^2(x) - y f_{\theta}(x)$$

$$(\nabla_{\theta} \phi)^T G \nabla_{\theta} \phi$$



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

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

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


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

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

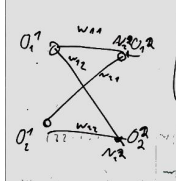
Norm von  $G$  ist unabhängig von  $G$

$$(t-y)^2$$

$$t-y$$

Proof of 3.4

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

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


$$\sigma'(N_1) N_1$$

$$\sigma'(N_2) N_2$$

$$\sigma^2 = N^2 \text{diag}(\sigma'(N^2))$$

$$\sigma^2 = \begin{bmatrix} \sigma^2_1 \\ \sigma^2_2 \end{bmatrix} = \begin{bmatrix} \sigma(N_1) \\ \sigma(N_2) \end{bmatrix} = \begin{bmatrix} \sigma'(N_1) N_1 \\ \sigma'(N_2) N_2 \end{bmatrix} = \begin{bmatrix} N_1^2 & N_2^2 \end{bmatrix}$$

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

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

Q.E.D.

[Proof of 3.5]

$x_i^T \in \mathbb{R}^p$

Netzwerk hat skalaren output

$X^T = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times p}$

$V \in \mathbb{R}^{p \times p}$

$X^T V \vec{0} = \vec{0}$

$\begin{bmatrix} x_1^T V \\ \vdots \\ x_n^T V \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

$\Rightarrow \sqrt{X} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \Rightarrow V^T X X^T V = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}$

[Proof of 3.6]

$(\log f)' = \frac{1}{f} \cdot f' = \frac{f'}{f} \quad (-1)$

$f = p_{\theta + \epsilon \alpha} \Big|_{t=0} \quad \left( \frac{d}{dt} \right)_{t=0} \text{ es gilt: } \frac{d p_{\theta + \epsilon \alpha}}{d \epsilon} \Big|_{\epsilon=0} = \bar{\alpha}$

$\Rightarrow \frac{d p_{\theta + \epsilon \alpha}}{d \epsilon} \Big|_{\epsilon=0} = \frac{\bar{\alpha}}{p_{\theta}} = \frac{f'}{f} = \frac{d}{d \epsilon} (\log p_{\theta + \epsilon \alpha}) \Big|_{\epsilon=0}$

$\Rightarrow \langle \bar{\alpha}, \bar{\beta} \rangle_{p_{\theta}} = \int \frac{\bar{\alpha}}{p_{\theta}} \frac{\bar{\beta}}{p_{\theta}} p_{\theta} = \int \frac{d}{d \epsilon} \log p_{\theta + \epsilon \alpha} \Big|_{\epsilon=0} \frac{d}{d \epsilon} \log p_{\theta + \epsilon \beta} \Big|_{\epsilon=0} p_{\theta}$

$M = \frac{1}{p_{\theta}}$

$E_{p_{\theta}}[\log p_{\theta + \epsilon \alpha} \log p_{\theta + \epsilon \beta}] = \int \log p_{\theta + \epsilon \alpha} \log p_{\theta + \epsilon \beta} p_{\theta}$

$\text{Abbildung}$