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

Sandro Braun, Leander Kurscheidt

July 17, 2018

**Abstract**

## 1 Geometry of Deep Rectified Networks

Lemma 2.1 (Structure in Gradient)

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

## References

$t = 1, \quad s = 1 \qquad 0 \le t \le s \le L$

$$O^{s+2}_{(x)} = \sigma\left( O_1^1 \, m_1 + O_2^1 \, m_2 \right)$$

$z$

$w_1$

$w_2$

$O^s$

$O^{s+1}$

$$\frac{\partial O^2}{\partial m_1} = \sigma'(z) \cdot O_1^1 \, w_1$$

$$\frac{\partial O^1}{\partial m_2} = \sigma'(z) \cdot O_2^1 \, w_2$$

$$= \sigma'(z) \left( O_1^1 \cdot m_1 + O_1^1 \, m_2 \right)$$

$z$

$$= \sigma(z)$$

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

$$\frac{\partial t}{\partial \theta} = \sigma'(x \cdot \theta) \cdot x \qquad \sigma(x \cdot \theta)$$

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

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

$$= f_\theta(x)$$

---

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

layer $t$

$N^t(x)$

$O^t(x)$

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

non-lin

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

$O^1$

$O^0$

$$\sigma\left( \sigma\left( X^T w^0 \right) w^1 \right)$$

$N^1$

$N^2$

$X \in \mathbb{R}^{28^2} \qquad Y \in \mathbb{R}^{10}$

$0 \quad 1 \quad 2$

$w^0 \, 628^2, \, w^1 \, 1024, \, w_2 = 10$

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

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

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

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

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

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

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

---

P.5 | Proof of Corollar 2.1

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

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

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

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

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

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

| Proof of Corollary 2.1

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

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

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

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

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

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

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

$\Updownarrow$

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

---

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

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

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

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

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

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

$m_{11}\theta_1^2 + \theta_2^2 m_{22}$

$+ (m_{12}+m_{21})\theta_1$
$\theta_2$

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

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

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

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

$$m_{11}\theta_1^2 + G_2^2 m_{22}$$
$$+ (m_{12}+m_{21})\theta_1$$

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

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

(Chain)
$$= \langle \frac{\partial l(f_\theta(X), Y)}{\partial f_\theta(X)} \nabla_\theta f_\theta(X), \theta \rangle$$

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

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

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

$s$

---

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

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

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

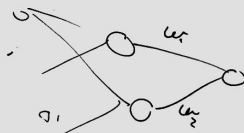$$l = \left(\frac{\ }{2}\right)^2 \to \frac{\partial l}{\partial f} = \mid \ \mid = \theta \ E\left[ \langle f_\theta(X)-Y, f_\theta(X) \rangle^2 \right]$$

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

$$\nabla_\theta \theta^{L+1} = \begin{bmatrix} \frac{\partial \theta_{L+1}}{\partial_{-1}} \\ \vdots \\ \partial u_n \end{bmatrix} \qquad \nabla\theta^T \theta = \qquad (\nabla\theta)^T G^T \theta$$

$s$

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

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

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



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

$$\vdots$$

$$= (L+1)^2 E\left( \left\langle \frac{\partial \ell(f(x_i, y_i))}{\partial f_\theta(x_i)}, f_\theta(x) \right\rangle^2 \right)$$

Norm
von
$\theta$

ist unabhängig von $\theta$

$$(f - y)^2$$

$$f - y$$