# Relative Fisher Information and Natural Gradient for Learning Large Modular Models

Ke Sun [1]    Frank Nielsen [2,3]

[1]King Abdullah University of Science & Technology (KAUST)

[2]École Polytechnique

[3]Sony CSL

ICML 2017

# Fisher Information Metric (FIM)

Consider a statistical model $p(\boldsymbol{x} \mid \boldsymbol{\Theta})$ of order $D$. The FIM (Hotelling29,Rao45) $\mathcal{I}(\boldsymbol{\Theta}) = (\mathcal{I}_{ij})$ is defined by a $D \times D$ positive semi-definite matrix

$$\mathcal{I}_{ij} = E_p \left[ \frac{\partial l}{\partial \Theta_i} \frac{\partial l}{\partial \Theta_j} \right], \tag{1}$$

where $l(\boldsymbol{\Theta}) = \log p(\boldsymbol{x} \mid \boldsymbol{\Theta})$ denotes the log-likelihood.

## Equivalent Expressions

$$\mathcal{I}_{ij} = E_p \left[ \frac{\partial l}{\partial \Theta_i} \frac{\partial l}{\partial \Theta_j} \right]$$

$$= -E_p \left[ \frac{\partial^2 l}{\partial \Theta_i \partial \Theta_j} \right]$$

$$= 4 \int \frac{\partial \sqrt{p(\boldsymbol{x} \mid \boldsymbol{\Theta})}}{\partial \Theta_i} \frac{\partial \sqrt{p(\boldsymbol{x} \mid \boldsymbol{\Theta})}}{\partial \Theta_j} d\boldsymbol{x}.$$

**Observed FIM (Efron & Hinkley, 1978)** With respect to $X_n = \{\boldsymbol{x}_k\}_{k=1}^n$,

$$\hat{\mathcal{I}} = -\nabla^2 l(\boldsymbol{\Theta} \mid X_n) = -\sum_{i=1}^n \frac{\partial^2 \log p(\boldsymbol{x}_i \mid \boldsymbol{\Theta})}{\partial \boldsymbol{\Theta} \partial \boldsymbol{\Theta}^{\mathsf{T}}}.$$

# FIM and Statistical Learning

▶ Any parametric learning is inside a corresponding parameter manifold $\mathcal{M}_\Theta$



$\mathcal{M}_\Theta$

$\mathcal{T}_{\boldsymbol{\theta}}\mathcal{M}_\Theta$: a tangent space with a local inner product $g(\boldsymbol{\theta})$

$\boldsymbol{\theta}$

a learning curve

▶ FIM gives an invariant Riemannian metric $g(\boldsymbol{\Theta}) = \mathcal{I}(\boldsymbol{\Theta})$ for any loss function based on standard f-divergence (KL, cross-entropy, . . . )

S. Amari. Information Geometry and Its Applications. 2016.

## Invariance

The FIM is *not* invariant and depends on the parameterization:

$$g_{\boldsymbol{\Theta}}(\boldsymbol{\Theta}) = \boldsymbol{J}^{\mathsf{T}} g_{\boldsymbol{\Lambda}}(\boldsymbol{\Lambda}) \boldsymbol{J}$$

where $\boldsymbol{J}$ is the Jacobian matrix $J_{ij} = \frac{\partial \Lambda_i}{\partial \Theta_j}$.

However its measurements such as $\langle \delta\boldsymbol{\Theta}, \delta\boldsymbol{\Theta} \rangle_{g(\boldsymbol{\Theta})}$ is invariant:

$$
\begin{aligned}
\langle \delta\boldsymbol{\Theta}, \delta\boldsymbol{\Theta} \rangle_{g(\boldsymbol{\Theta})} &= \delta\boldsymbol{\Theta}^{\mathsf{T}} g(\boldsymbol{\Theta}) \delta\boldsymbol{\Theta} \\
&= \delta\boldsymbol{\Theta}^{\mathsf{T}} \boldsymbol{J}^{\mathsf{T}} g_{\boldsymbol{\Lambda}}(\boldsymbol{\Lambda}) \boldsymbol{J} \delta\boldsymbol{\Theta} \\
&= \delta\boldsymbol{\Lambda}^{\mathsf{T}} g_{\boldsymbol{\Lambda}}(\boldsymbol{\Lambda}) \delta\boldsymbol{\Lambda} \\
&= \langle \delta\boldsymbol{\Lambda}, \delta\boldsymbol{\Lambda} \rangle_{g(\boldsymbol{\Lambda})}.
\end{aligned}
$$

Regardless of the choice of the coordinate system, it is essentially the same metric!

# Statistical Formulation of a Multilayer Perceptron (MLP)

$$p(\mathbf{y} \mid \mathbf{x}, \Theta) = \sum_{\mathbf{h}_1, \cdots, \mathbf{h}_{L-1}} p(\mathbf{y} \mid \mathbf{h}_{L-1}, \boldsymbol{\theta}_L) \cdots p(\mathbf{h}_2 \mid \mathbf{h}_1, \boldsymbol{\theta}_2) p(\mathbf{h}_1 \mid \mathbf{x}, \boldsymbol{\theta}_1),$$

# The FIM of a MLP

The FIM of a MLP has the following expression

$$g(\boldsymbol{\Theta}) = E_{\boldsymbol{x} \sim \hat{p}(X_n), \boldsymbol{y} \sim p(\boldsymbol{y} \mid \boldsymbol{x}, \boldsymbol{\Theta})} \left[ \frac{\partial l}{\partial \boldsymbol{\Theta}} \frac{\partial l}{\partial \boldsymbol{\Theta}^\mathsf{T}} \right]$$

$$= \frac{1}{n} \sum_{i=1}^{n} E_{p(\boldsymbol{y} \mid \boldsymbol{x}_i, \boldsymbol{\Theta})} \left[ \frac{\partial l_i}{\partial \boldsymbol{\Theta}} \frac{\partial l_i}{\partial \boldsymbol{\Theta}^\mathsf{T}} \right]$$

where

- $\hat{p}(X_n)$ is the empirical distribution of the samples $X_n = \{\boldsymbol{x}_i\}_{i=1}^{n}$
- $l_i(\boldsymbol{\Theta}) = \log p(\boldsymbol{y} \mid \boldsymbol{x}_i, \boldsymbol{\Theta})$ is the conditional log-likelihood

# Meaning of the FIM of a MLP

Consider a learning step on $\mathcal{M}_\Theta$ from $\Theta$ to $\Theta + \delta\Theta$. The step size

$$
\begin{aligned}
\langle \delta\Theta, \delta\Theta \rangle_{g(\Theta)} &= \delta\Theta^\intercal g(\Theta)\delta\Theta \\
&= \delta\Theta^\intercal \left\{ \frac{1}{n} \sum_{i=1}^{n} E_{p(y \mid x_i, \Theta)} \left[ \frac{\partial l_i}{\partial \Theta} \frac{\partial l_i}{\partial \Theta^\intercal} \right] \right\} \delta\Theta \\
&= \frac{1}{n} \sum_{i=1}^{n} E_{p(y \mid x_i, \Theta)} \left[ \delta\Theta^\intercal \frac{\partial l_i}{\partial \Theta} \right]^2
\end{aligned}
$$

measures how much $\delta\Theta$ is statistically along $\frac{\partial l}{\partial \Theta}$.

**Will $\delta\Theta$ make a significant change to the mapping $x \rightarrow y$ or not?**

# Natural Gradient: Seeking a Short Path

Consider $\min_{\Theta \in \mathcal{M}_\Theta} L(\Theta)$. At $\Theta_t \in \mathcal{M}_\Theta$, the target is to minimize wrt $\delta\Theta$

$$\underbrace{L(\Theta_t + \delta\Theta)}_{\text{Loss function}} + \frac{1}{2\gamma} \underbrace{\langle \delta\Theta, \delta\Theta \rangle_{g(\Theta_t)}}_{\text{Squared step size}} \quad (\gamma\text{: learning rate})$$

$$\approx L(\Theta_t) + \delta\Theta^\mathsf{T} \bigtriangledown L(\Theta_t) + \frac{1}{2\gamma} \delta\Theta^\mathsf{T} g(\Theta_t) \delta\Theta,$$

giving a learning step

$$\delta\Theta_t = -\gamma \underbrace{g^{-1}(\Theta_t) \bigtriangledown L(\Theta_t)}_{\text{natural gradient}}$$

▶ Equivalence with mirror descent (Raskutti & Mukherjee 2013)

# Natural Gradient: Intrinsics

$$\delta\boldsymbol{\Theta}_t = -\gamma g^{-1}(\boldsymbol{\Theta}_t) \bigtriangledown L(\boldsymbol{\Theta}_t)$$

This Riemannnian metric is a property of the parameter space that is independent of the loss function $L(\boldsymbol{\Theta})$.

The good performance of natural gradient relies on that $L(\boldsymbol{\Theta})$ is similarly curved as $\log p(\boldsymbol{x} \,|\, \boldsymbol{\Theta})$ ($\boldsymbol{x} \sim p(\boldsymbol{x} \,|\, \boldsymbol{\Theta})$).

Natural gradient is not universally good for any loss functions.

# Natural Gradient: Pros and Cons

### Pros

- Invariant (intrinsic) gradient
- Not trapped in plateaus
- Achieve Fisher efficiency in online learning

### Cons

- Too expensive to compute (no closed-form FIM; need matrix inversion)

# Relative FIM — Informal Ideas

- Decompose the learning system into subsystems

- The subsystems are interfaced with each other through hidden variables $\boldsymbol{h}_i$

- Some subsystems are interfaced with the I/O environment through $\boldsymbol{x}_i$ and $\boldsymbol{y}_i$

- Compute the subsystem FIM by **integrating out its interface variables $\boldsymbol{h}_i$**, so that the intrinsics of this subsystem can be discussed regardless of the remaining parts

# From FIM to Relative FIM (RFIM)

### FIM

$\boldsymbol{\theta}$ $\longrightarrow$ $\triangle$ $\longrightarrow$ $\log p(\boldsymbol{x} \,|\, \boldsymbol{\theta})$ (likelihood scalar)

(parameter vector)

How sensitive is $\boldsymbol{x}$ wrt tiny movements of $\boldsymbol{\theta}$ on $\mathcal{M}_{\boldsymbol{\theta}}$?

### RFIM

$\boldsymbol{\theta}$ $\longrightarrow$ $\triangle$ $\longrightarrow$ $\log p(\boldsymbol{r} \,|\, \boldsymbol{\theta}, \boldsymbol{\theta}_f)$ (likelihood scalar)

(parameter vector)

Given $\boldsymbol{\theta}_f$, how sensitive is $\boldsymbol{r}$ wrt tiny movements of $\boldsymbol{\theta}$?

## Relative FIM — Definition

Given $\theta_f$ (the **reference**), the Relative Fisher Information Metric (RFIM) of $\theta$ wrt $h$ (the **response**) is

$$g^h(\theta \,|\, \theta_f) = E_{p(h \,|\, \theta, \theta_f)}\left[\frac{\partial}{\partial \theta} \ln p(h \,|\, \theta, \theta_f)\frac{\partial}{\partial \theta^\intercal} \ln p(h \,|\, \theta, \theta_f)\right],$$

or simply $g^h(\theta)$.

Meaning: given $\theta_f$, how variations of $\theta$ will affect the response $h$.

# Different Subsystems – Simple Examples



**Figure:** Generator

**Figure:** Discriminator or Regressor

# A Dynamic Geometry



Model: $p(\boldsymbol{y} \mid \boldsymbol{\Theta}, \boldsymbol{x}) = \sum_{\boldsymbol{h}_1} \sum_{\boldsymbol{h}_2} p(\boldsymbol{h}_1 \mid \boldsymbol{\theta}_1, \boldsymbol{x}) \quad p(\boldsymbol{h}_2 \mid \boldsymbol{\theta}_2, \boldsymbol{h}_1) \quad p(\boldsymbol{y} \mid \boldsymbol{\theta}_3, \boldsymbol{h}_2)$

- As the interface hidden variables $\boldsymbol{h}_i$ are changing, the subsystem geometry is not absolute but is **relative** to its reference variables provided by adjacent subsystems

# RFIM of One `tanh` Neuron

Consider a neuron with input $\boldsymbol{x}$, weights $\boldsymbol{w}$, a hyperbolic tangent activation function, and a stochastic output $y \in \{-1, 1\}$, given by

$$p(y = 1) = \frac{1 + \tanh(\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}})}{2}, \quad \tanh(t) = \frac{\exp(t) - \exp(-t)}{\exp(t) + \exp(-t)}.$$

$\tilde{\boldsymbol{x}} = (\boldsymbol{x}^\mathsf{T}, 1)^\mathsf{T}$ denotes the augmented vector of $\boldsymbol{x}$

$$g^y(\boldsymbol{w} \,|\, \boldsymbol{x}) = \nu_{\tanh}(\boldsymbol{w}, \boldsymbol{x})\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T}, \quad \nu_{\tanh}(\boldsymbol{w}, \boldsymbol{x}) = \operatorname{sech}^2(\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}}).$$

# RFIM of Parametric Rectified Linear Unit

$$p(y \mid \boldsymbol{w}, \boldsymbol{x}) = G(y \mid \texttt{relu}(\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}}), \sigma^2), \quad (G \text{ is for Gaussian})$$

$$\texttt{relu}(t) = \begin{cases} t & \text{if } t \geq 0 \\ \iota t & \text{if } t < 0. \end{cases} \quad (0 \leq \iota < 1)$$

By certain assumptions,

$$g^y(\boldsymbol{w} \mid \boldsymbol{x}) = \nu_{\texttt{relu}}(\boldsymbol{w}, \boldsymbol{x})\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T},$$

$$\nu_{\texttt{relu}}(\boldsymbol{w}, \boldsymbol{x}) = \frac{1}{\sigma^2}\left[\iota + (1 - \iota)\underbrace{\texttt{sigm}}_{sigmoid}\left(\frac{1-\iota}{\omega}\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}}\right)\right]^2.$$

Set $\sigma = 1$, $\iota = 0$, it simplifies to

$$\nu_{\texttt{relu}}(\boldsymbol{w}, \boldsymbol{x}) = \texttt{sigm}^2\left(\frac{1}{\omega}\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}}\right).$$

## Generic Expression of One-neuron RFIMs

Denote $f \in \{\mathtt{tanh}, \mathtt{sigm}, \mathtt{relu}, \mathtt{elu}\}$ to be an element-wise nonlinear activation function. The RFIM is

$$g^y(\mathbf{w} \mid \mathbf{x}) = \nu_f(\mathbf{w}, \mathbf{x})\tilde{\mathbf{x}}\tilde{\mathbf{x}}^\mathsf{T},$$

where $\nu_f(\mathbf{w}, \mathbf{x})$ is a positive coefficient with large values in the *linear region*, or the effective learning zone of the neuron.

## RFIM of a Linear Layer

$x$: input; $W$: connection weights; $y$: stochastic output following

$$p(y \mid W, x) = G(y \mid W^\mathsf{T}\tilde{x}, \sigma^2 I).$$

We vectorize $W$ by stacking its columns $\{w_i\}$. Then

$$g^{y}(W \mid x) = \frac{1}{\sigma^2} \begin{bmatrix} \tilde{x}\tilde{x}^\mathsf{T} & & \\ & \ddots & \\ & & \tilde{x}\tilde{x}^\mathsf{T} \end{bmatrix}.$$

# RFIM of a Non-linear Layer

A *nonlinear* layer applies an element-wise activation on $\boldsymbol{W}^\intercal \tilde{\boldsymbol{x}}$. We have

$$g^{\boldsymbol{y}}\left(\boldsymbol{W} \mid \boldsymbol{x}\right) = \begin{bmatrix} \nu_f(\boldsymbol{w}_1, \boldsymbol{x})\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\intercal & & \\ & \ddots & \\ & & \nu_f(\boldsymbol{w}_m, \boldsymbol{x})\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\intercal \end{bmatrix},$$

where $\nu_f(\boldsymbol{w}_i, \boldsymbol{x})$ depends on the activation function $f$.

The RFIMs of single neuron models, a linear layer, a non-linear layer, a soft-max layer, two consecutive layers all have **simple closed form solutions**[1].

# List of RFIMs

| Subsystem | the RFIM $g^{\boldsymbol{y}}(\boldsymbol{w})$ |
|---|---|
| A tanh neuron | $\operatorname{sech}^2(\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}})\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T}$ |
| A sigm neuron | $\operatorname{sigm}(\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}})\left[1 - \operatorname{sigm}(\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}})\right]\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T}$ |
| A relu neuron | $\left[\iota + (1-\iota)\operatorname{sigm}\left(\frac{1-\iota}{\omega}\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}}\right)\right]^2\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T}$ |
| A elu neuron | $\begin{cases} \tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} & \text{if } \boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}} \geq 0 \\ (\alpha\exp(\boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}}))^2\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} & \text{if } \boldsymbol{w}^\mathsf{T}\tilde{\boldsymbol{x}} < 0 \end{cases}$ |
| A linear layer | $\operatorname{diag}\left[\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T}, \cdots, \tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T}\right]$ |
| A non-linear layer | $\operatorname{diag}\left[\nu_f(\boldsymbol{w}_1, \tilde{\boldsymbol{x}})\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T}, \cdots, \nu_f(\boldsymbol{w}_m, \tilde{\boldsymbol{x}})\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T}\right]$ |
| A soft-max layer | $\begin{bmatrix} (\eta_1 - \eta_1^2)\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} & -\eta_1\eta_2\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} & \cdots & -\eta_1\eta_m\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} \\ -\eta_2\eta_1\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} & (\eta_2 - \eta_2^2)\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} & \cdots & -\eta_2\eta_m\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} \\ \vdots & \vdots & \ddots & \vdots \\ -\eta_m\eta_1\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} & -\eta_m\eta_2\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} & \cdots & (\eta_m - \eta_m^2)\tilde{\boldsymbol{x}}\tilde{\boldsymbol{x}}^\mathsf{T} \end{bmatrix}.$ |
| Two layers | see the paper. |

# Relative Natural Gradient Descent (RNGD)

For each subsystem,

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \gamma \cdot \underbrace{\left(\bar{g}^{\boldsymbol{h}}(\boldsymbol{\theta}_t \,|\, \boldsymbol{\theta}_f)\right)^{-1}}_{\text{inverse RFIM}} \cdot \left.\frac{\partial L}{\partial \boldsymbol{\theta}}\right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}$$

where

$$\bar{g}^{\boldsymbol{h}}(\boldsymbol{\theta}_t \,|\, \boldsymbol{\theta}_f) = \frac{1}{n}\sum_{i=1}^{n} g^{\boldsymbol{h}}(\boldsymbol{\theta}_t \,|\, \boldsymbol{\theta}_f^i).$$

By definition, RFIM is a function of the reference variables.
$\bar{g}^{\boldsymbol{h}}(\boldsymbol{\theta}_t \,|\, \boldsymbol{\theta}_f)$ is its expectation wrt an empirical distribution of $\boldsymbol{\theta}_f$.

# Proof-of-concept



- ▶ MLP with shape 784-80-80-80-10
- ▶ `relu` activation
- ▶ Mini batch size 50
- ▶ Recompute the inverse RFIM every 100 mini batchs
- ▶ $L_2$ regularization

# BNA: batch normalization (BN) after activation

# Change the MLP shape to 784-100-100-100-10

# Novel Viewpoint

Learning is a process where a set of collaborative learners move on their sub-manifolds, and the geometries of these sub-manifolds are also evolving with the system.

- ▶ Well-suited to parallel computation and distributed learning

# Conclusion

- FIM is just a special case of RFIM, where the subsystem is the whole system

- By looking at smaller subsystems, RFIM can have simpler closed-form expressions

- RNGD can be implemented without approximation

- This has the potential to improve learning of large neural networks

codes, updates:

https://www.lix.polytechnique.fr/~nielsen/RFIM/

**Thank you!**