

---

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

---

Ke Sun<sup>1</sup> Frank Nielsen<sup>2,3</sup>

## Abstract

Fisher information and natural gradient provided deep insights and powerful tools to artificial neural networks. However related analysis becomes more and more difficult as the learner’s structure turns large and complex. This paper makes a preliminary step towards a new direction. We extract a local component from a large neural system, and define its relative Fisher information metric that describes accurately this small component, and is invariant to the other parts of the system. This concept is important because the geometry structure is much simplified and it can be easily applied to guide the learning of neural networks. We provide an analysis on a list of commonly used components, and demonstrate how to use this concept to further improve optimization.

## 1. Fisher Information Metric

The Fisher Information Metric (FIM)  $\mathcal{I}(\Theta) = (\mathcal{I}_{ij})$  of a statistical parametric model  $p(\mathbf{x} | \Theta)$  of order  $D$  is defined by a  $D \times D$  positive semidefinite (psd) matrix ( $\mathcal{I}(\Theta) \succeq 0$ ) with coefficients  $\mathcal{I}_{ij} = E_p \left[ \frac{\partial l}{\partial \Theta_i} \frac{\partial l}{\partial \Theta_j} \right]$ , where  $l(\Theta)$  denotes the log-density function  $\log p(\mathbf{x} | \Theta)$ . Under light regularity conditions, FIM can be rewritten equivalently as

$$\mathcal{I}_{ij} = -E_p \left[ \frac{\partial^2 l}{\partial \Theta_i \partial \Theta_j} \right] = 4 \int \frac{\partial \sqrt{p(\mathbf{x} | \Theta)}}{\partial \Theta_i} \frac{\partial \sqrt{p(\mathbf{x} | \Theta)}}{\partial \Theta_j} d\mathbf{x}.$$

As its empirical counterpart, the observed FIM (Efron & Hinkley, 1978) with respect to (wrt) a sample set  $X_n = \{\mathbf{x}_k\}_{k=1}^n$  is  $\hat{\mathcal{I}}(\Theta | X_n) = -\nabla^2 l(\Theta | X_n)$ , which is often evaluated at the maximum likelihood estimate  $\hat{\Theta} = \hat{\Theta}(X_n)$ . By the law of large numbers,  $\hat{\mathcal{I}}(\Theta)$  converges to the (expected) FIM  $\mathcal{I}(\Theta)$  as  $n \rightarrow \infty$ .

---

<sup>1</sup>King Abdullah University of Science and Technology (KAUST), Saudi Arabia <sup>2</sup>École Polytechnique, France <sup>3</sup>Sony Computer Science Laboratories Inc., Japan. Correspondence to: Ke Sun <sunk@ieec.org>, Frank Nielsen <Frank.Nielsen@acm.org>.

The FIM is *not* invariant and depends on the parameterization. We can optionally write  $\mathcal{I}(\Theta)$  as  $\mathcal{I}_{\Theta}(\Theta)$  to emphasize the coordinate system. By definition,  $\mathcal{I}_{\Theta}(\Theta) = \mathbf{J}^T \mathcal{I}_{\Lambda}(\Lambda) \mathbf{J}$  where  $\mathbf{J} = (J_{ij})$ ,  $J_{ij} = \frac{\partial \Lambda_i}{\partial \Theta_j}$  is the Jacobian matrix. For example, the FIM of regular natural exponential families (NEFs)  $l(\Theta) = \Theta^T t(\mathbf{x}) - F(\Theta)$  (log-linear models with sufficient statistics  $t(\mathbf{x})$ ) is  $\mathcal{I}(\Theta) = \nabla^2 F(\Theta) \succ 0$ , the Hessian of the log-normalizer function  $F(\Theta)$ . Although exponential families can approximate arbitrarily *any* smooth density (Cobb et al., 1983), the log-normalizer function may not be available in closed-form nor computationally tractable (Montanari, 2015).

The FIM is an important concept for statistical machine learning. It gives a Riemannian metric (Hotelling, 1929; Rao, 1945) of the learning parameter space which is unique (Čencov, 1982; Dowty, 2017). Hence *any* learning is in a space that is intrinsically curved based on the FIM, regardless of the choice of the coordinate system. It also gives a bound (Fréchet, 1943; Cramér, 1946; Nielsen, 2013) of learning efficiency saying that the variance of *any* unbiased learning of  $\Theta$  is at least  $\mathcal{I}^{-1}(\Theta)/n$ , where  $n$  is the i.i.d. sample size. The FIM is applied to neural network optimization (Amari, 1997), metric learning (Lebanon, 2005), reinforcement learning (Thomas, 2014) and manifold learning (Sun & Marchand-Maillet, 2014).

However computing the FIM is expensive. Besides the fact that learning machines have often singularities (Watanabe, 2009) ( $|\mathcal{I}(\Theta)| = 0$ , not full rank) characterized by plateaux in gradient learning, computing/estimating the FIM of a large neuron system (e.g. one with millions of parameters, Szegedy, Christian et al. 2015) is very challenging due to the finiteness of data, and the huge number  $\frac{D(D+1)}{2}$  of matrix coefficients to evaluate. Furthermore, gradient descent techniques require inverting this large matrix and tuning the learning rate.

To tackle this problem, past works mainly focus on how to approximate the FIM with a block diagonal form (Kurita, 1994; Le Roux et al., 2008; Martens, 2010; Pascanu & Bengio, 2014; Martens & Grosse, 2015) or quasi-diagonal form (Ollivier, 2013; Marceau-Caron & Ollivier, 2016). This global approach faces increasing approximation error and increasing computational cost as the system scales up

and as complex and dynamic structures (Looks et al., 2017) emerge.

This work aims at a different local approach. The idea is to accurately describe the information geometry (IG) in a subsystem of the large learning system, which is invariant to the scaling up and structural change of the global system, so that the local machinery, including optimization, can be discussed regardless of the other parts.

For this purpose, a novel concept, the Relative Fisher Information Metric (RFIM), is defined. Unlike the traditional geometric view of a high-dimensional parameter manifold, RFIMs defines *multiple projected low-dimensional geometries of subsystems*. This geometry is correlated to the parameters beyond the subsystem and is therefore considered *dynamic*. It can be used to characterize the efficiency of a local learning process. Taking this stance has potential in deep learning because a deep neural network can be decomposed into many local components such as neurons or layers. The RFIM is well suited to the compositional block structures of neural networks. The RFIM can be used for out-of-core learning.

The paper is organized as follows. Sec. 2 reviews natural gradient within the context of Multi-Layer Perceptrons (MLPs). Sec. 3 formally defines the RFIM, and gives a table of RFIMs of several commonly used subsystems. Sec. 4 discusses the advantages of using the RFIM as compared to the FIM. Sec. 5 gives an algorithmic framework and proof-of-concept experiments on neural network optimization. Sec. 6 presents related works on parameter diagonalization. Sec. 7 concludes this work and further hints at perspectives.

## 2. Natural Gradient: Review and Insights

Consider a MLP  $\mathbf{x} \xrightarrow{\theta_1} \mathbf{h}_1 \cdots \mathbf{h}_{L-1} \xrightarrow{\theta_L} \mathbf{y}$ , whose statistical model is the following conditional distribution

$$p(\mathbf{y} | \mathbf{x}, \Theta) = \sum_{\mathbf{h}_1, \dots, \mathbf{h}_{L-1}} p(\mathbf{h}_1 | \mathbf{x}, \theta_1) \cdots p(\mathbf{y} | \mathbf{h}_{L-1}, \theta_L).$$

The often intractable sum over  $\mathbf{h}_1, \dots, \mathbf{h}_{L-1}$  can be get rid off by deteriorating  $p(\mathbf{h}_1 | \mathbf{x}, \theta_1), \dots, p(\mathbf{h}_{L-1} | \mathbf{h}_{L-2}, \theta_{L-1})$  to Dirac's deltas  $\delta$ , and letting merely the last layer  $p(\mathbf{y} | \mathbf{h}_{L-1}, \theta_L)$  be stochastic. Other models such as restricted Boltzmann machines (Nair & Hinton, 2010; Montavon & Müller, 2012), deep belief networks (Hinton et al., 2006), dropout (Wager et al., 2013), and variational autoencoders (Kingma & Welling, 2014) do consider the  $\mathbf{h}_i$ 's to be stochastic.

The tensor metric of the *neuromanifold* (Amari, 1995)  $\mathcal{M}$ , consisting of all MLPs with the same architecture but different parameter values, is locally defined by the FIM. Because a MLP corresponds to a con-

ditional distribution, its FIM is a function of the input  $\mathbf{x}$ . By taking an empirical average over the input samples  $\{\mathbf{x}_k\}_{k=1}^n$ , the FIM of a MLP can be expressed as  $\mathcal{I}_{\Theta}(\Theta) = \frac{1}{n} \sum_{k=1}^n E_{p(\mathbf{y} | \mathbf{x}_k, \Theta)} \left[ \frac{\partial l_k}{\partial \Theta} \frac{\partial l_k}{\partial \Theta^\top} \right]$ , where  $l_k(\Theta) = \log p(\mathbf{y} | \mathbf{x}_k, \Theta)$  denotes the conditional log-likelihood function wrt  $\mathbf{x}_k$ .

To understand the meaning of the Riemannian metric  $\mathcal{I}_{\Theta}(\Theta)$ , it measures the intrinsic difference between two nearby neural networks around  $\Theta \in \mathcal{M}$ . A learning step can be regarded as a tiny displacement  $\delta\Theta$  on  $\mathcal{M}$ . According to the FIM, the infinitesimal square distance

$$\langle \delta\Theta, \delta\Theta \rangle_{\mathcal{I}_{\Theta}(\Theta)} = \frac{1}{n} \sum_{k=1}^n E_{p(\mathbf{y} | \mathbf{x}_k, \Theta)} \left[ \left( \delta\Theta^\top \frac{\partial l_k}{\partial \Theta} \right)^2 \right] \quad (1)$$

measures how much  $\delta\Theta$  (with a radius constraint) is statistically along  $\frac{\partial l}{\partial \Theta}$ , or equivalently how much  $\delta\Theta$  affects intrinsically the conditional distribution  $p(\mathbf{y} | \mathbf{x}, \Theta)$ .

Consider the negative log-likelihood function  $L(\Theta) = -\sum_{k=1}^n \log p(\mathbf{y}_k | \mathbf{x}_k, \Theta)$  wrt the observed pairs  $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k=1}^n$ , we try to minimize the loss while maintaining a small learning step size  $\langle \delta\Theta, \delta\Theta \rangle_{\mathcal{I}_{\Theta}(\Theta)}$  on  $\mathcal{M}$ . At  $\Theta_t \in \mathcal{M}$ , the target is to minimize wrt  $\delta\Theta$  the Lagrange function

$$\begin{aligned} & L(\Theta_t + \delta\Theta) + \frac{1}{2\gamma} \langle \delta\Theta, \delta\Theta \rangle_{\mathcal{I}_{\Theta}(\Theta_t)} \\ & \approx L(\Theta_t) + \delta\Theta^\top \nabla_{\Theta} L(\Theta_t) + \frac{1}{2\gamma} \delta\Theta^\top \mathcal{I}_{\Theta}(\Theta_t) \delta\Theta, \end{aligned}$$

where  $\gamma > 0$  is a learning rate. The optimal solution of the above quadratic optimization gives a learning step

$$\delta\Theta_t = -\gamma \mathcal{I}_{\Theta}^{-1}(\Theta_t) \nabla_{\Theta} L(\Theta_t).$$

In this update procedure,  $\tilde{\nabla}_{\Theta} L(\Theta) = \mathcal{I}_{\Theta}^{-1}(\Theta) \nabla_{\Theta} L(\Theta)$  replaces the role of the usual gradient  $\nabla_{\Theta} L(\Theta)$  and is called the *natural gradient* (Amari, 1997).

Although the FIM depends on the chosen parameterization, the natural gradient is *invariant* to reparameterization. Let  $\Lambda$  be another coordinate system and  $\mathbf{J}$  be the Jacobian matrix of the mapping  $\Theta \rightarrow \Lambda$ . Then we have

$$\begin{aligned} \mathcal{I}_{\Theta}^{-1}(\Theta) \nabla_{\Theta} L(\Theta) &= (\mathbf{J}^\top \mathcal{I}_{\Lambda}(\Lambda) \mathbf{J})^{-1} \mathbf{J}^\top \nabla_{\Lambda} L(\Lambda) \\ &= \mathbf{J}^{-1} \mathcal{I}_{\Lambda}^{-1}(\Lambda) \nabla_{\Lambda} L(\Lambda), \end{aligned}$$

showing that  $\tilde{\nabla}_{\Theta} L(\Theta)$  and  $\tilde{\nabla}_{\Lambda} L(\Lambda)$  are the same dynamic up to coordinate transformation. As the learning rate  $\gamma$  is not infinitesimal in practice, natural gradient descent actually depends on the coordinate system (see e.g. Martens 2014). Other intriguing properties of natural gradient optimization lie in being free from getting trapped in plateaux of the error surface, and attaining Fisher efficiency in on-line learning (see Sec. 4 Amari 1998).

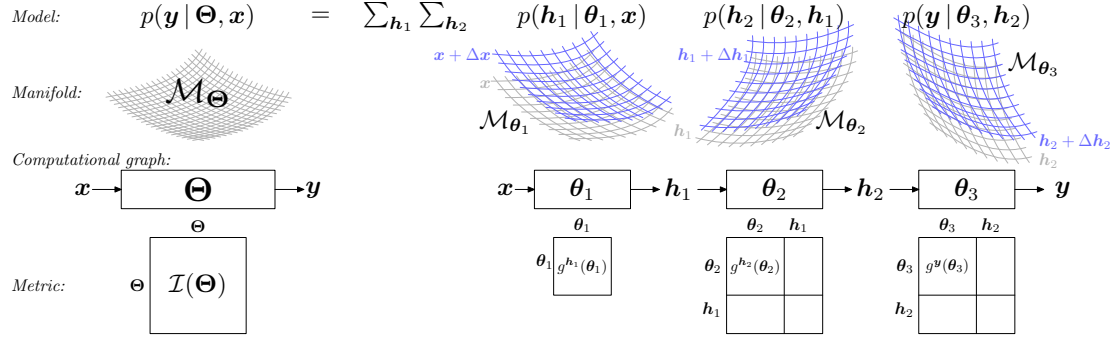


Figure 1. (left) The traditional global geometry of a MLP; (right) information geometry of subsystems. The gray and blue meshes show that the subsystem geometry is *dynamic* when the reference variable makes a tiny move. The square under the (sub-)system means the (R-)FIM is computed by (i) computing the FIM in the traditional way wrt all free parameters that affect the system output; (ii) choosing a sub-block that contains only the internal parameters of the (sub-)system and regarding the remaining variables as the reference.

For the sake of simplicity, we do not discuss singular FIMs with a subset of parameters having zero metric. This set of parameters forms an analytic variety (Watanabe, 2009), and technically the MLP as a statistical model is said to be non-regular (and the parameter  $\Theta$  is not identifiable). The natural gradient has been extended (Thomas, 2014) to cope with singular FIMs having positive semi-definite matrices by taking the Moore-Penrose pseudo-inverse (that coincides with the inverse matrix for full rank matrices).

In the family of 2<sup>nd</sup>-order optimization methods, a fuzzy line can be drawn from the natural gradient and alternative methods such as the Hessian-free optimization (Martens, 2010). By definition, the FIM is a property of the parameter space which is independent or weakly dependent on the input samples. For example, the FIM of a MLP is independent of  $\{y_i\}$ . In contrast, the Hessian (or related concepts such as the Gauss-Newton matrix, Martens 2014) is a property of the learning cost function wrt the input samples.

Bonnabel (Bonnabel, 2013) proposed to use the Riemannian exponential map to define a gradient descent step, thus ensuring to stay on the manifold for any chosen learning rate. Convergence is proven for Hadamard manifolds (of negative curvatures). However, it is not mathematically tractable to express the exponential map of hierarchical model manifolds like the neuromanifold.

### 3. RFIM: Definition and Expressions

In general, for large parametric systems, it is impossible to diagonalize or decorrelate all the parameters, so that we split instead all random variables into three parts  $\theta_f, \theta$  and  $h$ . We examine their intuitive meanings before giving the formal definition. The *reference*,  $\theta_f$ , consists of the majority of the random variables that are considered fixed (therefore allowing us to simplify the analysis). This is in analogy to the notion of a reference frame in physics.  $\theta$  is the

subsystem parameters, resembling the long-term memory adapting slowly to the observations (e.g. neural network weights). The *response*  $h$  is a random variable that reacts to the variations of  $\theta$ . Usually,  $h$  is the output of the subsystem that is connected to neighbour subsystems (e.g. hidden layer outputs). Formally, a subsystem which factorizes the learning machine is characterized by the conditional distribution  $p(h | \theta, \theta_f)$ , where  $\theta$  can be estimated based on  $h$  and  $\theta_f$ . We make the following definition.

**Definition 1 (RFIM).** Given  $\theta_f$ , the RFIM<sup>1</sup> of  $\theta$  wrt  $h$  is

$$g^h(\theta | \theta_f) \stackrel{\text{def}}{=} E_{p(h | \theta, \theta_f)} \left[ \frac{\partial}{\partial \theta} \log p(h | \theta, \theta_f) \frac{\partial}{\partial \theta^\top} \log p(h | \theta, \theta_f) \right],$$

or simply  $g^h(\theta)$ , corresponding to the estimation of  $\theta$  based on observations of  $h$  given  $\theta_f$ .

For example, consider a MLP. If we choose  $\theta_f$  to be the input features  $x$ , choose  $h$  to be the final output  $y$ , and choose  $\theta$  to be all the network weights  $\Theta$ , then the RFIM becomes the FIM:  $\mathcal{I}(\Theta) = g^y(\Theta | x)$ .

More generally, we can choose the response  $h$  to be other than the observables to compute the Fisher information of subsystems, especially dynamically during the learning of the global machine. To see the meaning of the RFIM, similar to eq. (1), the infinitesimal square distance  $\langle \delta\theta, \delta\theta \rangle_{g^h(\theta)} = E_{p(h | \theta, \theta_f)} \left[ (\delta\theta^\top \frac{\partial}{\partial \theta} \log p(h | \theta, \theta_f))^2 \right]$  measures how much  $\delta\theta$  impacts intrinsically the stochastic mapping  $\theta \rightarrow h$  which features the subsystem. We have the following proposition following definition 1.

**Proposition 2 (Relative Geometry Consistency).** If  $\theta_1$  consists of a subset of  $\theta_2$  so that  $\theta_2 = (\theta_1, \hat{\theta}_1)$ , then  $\forall \hat{\theta}_1, \mathcal{M}_{\theta_1}$  with the metric  $g^h(\theta_1 | \hat{\theta}_1)$  has exactly the same Rie-

<sup>1</sup>We use the same term ‘‘relative FIM’’ (Zegers, 2015) with a different definition.

mannian metric with the sub-manifold  $\{\theta_2 \in \mathcal{M}_{\theta_2} : \bar{\theta}_1 \text{ is fixed}\}$  induced by the ambient metric  $g^h(\theta_2)$ .

When the response  $h$  is chosen, then different splits of  $(\theta, \theta_f)$  are consistent with the same ambient geometry.

Figure 1 shows the traditional global geometry of a learning system, where the curvature is defined by the learner’s parameter sensitivity to the external environment ( $x$  and  $y$ ), as compared to the information geometry of subsystems, where the curvature is defined by the parameter sensitivity wrt hidden interface variables  $h$ . The two-colored meshes show that the geometry structure is dynamic and varies with the reference variable  $\theta_f$ .

One should not confuse the RFIM with the diagonal blocks of the FIM (Kurita, 1994). Both their meanings and expressions are different. The RFIM is computed by integrating out the hidden response variables  $h$ . The FIM is always computed by integrating out the observables  $x$  and  $y$ . Hence the RFIM is a more general concept and includes the FIM as a special case. This highlights a main difference with the backpropagated metric (Ollivier, 2013), which essentially considers parameter sensitivity wrt the final output. Despite the fact that the FIMs of small parametric structures such as single neurons was studied (Amari, 1997), we are not looking at a small single-component system but a component embedded in a large system, targeting at improving the large system.

In the following we provide a short table of commonly used RFIMs for future reference (the RFIMs listed are mostly straightforward from definition 1, with detailed derivations given in the supplementary material). This is meaningful since the RFIM is a new concept. We also want to demonstrate these simple closed form expressions without any approximations.

### 3.1. RFIMs of One Neuron

We start from the RFIM of single neuron models. Consider a stochastic neuron with input  $x$  and weights  $w$ . After a nonlinear activation function  $f$ , the output  $y$  is randomized surrounding the mean  $f(w^\top \tilde{x})$  with a variance. Throughout this paper  $\tilde{x} = (x^\top, 1)^\top$  denotes the augmented vector of  $x$  (homogeneous coordinates) so that  $w^\top \tilde{x}$  contains a bias term, and a general linear transformation can be written simply as  $A\tilde{x}$ .

Using  $x$  as the reference, the RFIM of  $w$  with respect to  $y$  has a common form  $g^y(w | x) = \nu_f(w, x) \tilde{x} \tilde{x}^\top$ , where  $\nu_f(w, x)$  is a positive coefficient with large values in the linear region, or the effective learning zone of the neuron. This agrees with early studies on single neuron FIMs (Amari, 1997; Kurita, 1994).

If  $f(t) = \tanh(t)$  is the hyperbolic tangent func-

tion, then  $\nu_f(w, x) = \text{sech}^2(w^\top \tilde{x})$ , where  $\text{sech}(t) = \frac{2}{\exp(t) + \exp(-t)}$  is the hyperbolic secant function. Similarly, if  $f(t) = \text{sigm}(t)$  is the sigmoid function, then  $\nu_f(w, x) = \text{sigm}(w^\top \tilde{x}) [1 - \text{sigm}(w^\top \tilde{x})]$ .

If  $f$  is defined by Parametric Rectified Linear Unit (PReLU) (He et al., 2015), which includes Rectified Linear Unit (ReLU) (Nair & Hinton, 2010) as a special case, so that  $f(t) = t (t \geq 0)$ ,  $f(t) = \iota t (t < 0)$ ,  $0 \leq \iota < 1$ , then under certain approximations (see supplementary material)

$$\nu_f(w, x) = \left[ \iota + (1 - \iota) \text{sigm} \left( \frac{1 - \iota}{\omega} w^\top \tilde{x} \right) \right]^2,$$

where  $\omega > 0$  is a hyper-parameter (e.g.  $\omega = 1$ ).

For the exponential linear unit (ELU) (Clevert et al., 2015),  $f(t) = t (t \geq 0)$ ,  $f(t) = \alpha (\exp(t) - 1) (t < 0)$ , where  $\alpha > 0$  is a hyper-parameter. We get

$$\nu_f(w, x) = \begin{cases} 1 & \text{if } w^\top \tilde{x} \geq 0 \\ \alpha^2 \exp(2w^\top \tilde{x}) & \text{if } w^\top \tilde{x} < 0. \end{cases}$$

### 3.2. RFIM of One Layer

Let  $D$  denote the dimensionality of the corresponding variable. A linear layer with input  $x$ , connection weights  $W = [w_1, \dots, w_{D_y}]$ , and stochastic output  $y$  can be represented by  $y \sim G(W^\top \tilde{x}, \sigma^2 I)$ , where  $I$  is the identity matrix, and  $\sigma$  is the scale of the observation noise, and  $G(\mu, \Sigma)$  is a multivariate Gaussian distribution with mean  $\mu$  and covariance matrix  $\Sigma$ . We vectorize  $W$  by stacking its columns  $\{w_i\}$ . Then  $g^y(W | x)$  is a tensor of size  $(D_x + 1)D_y \times (D_x + 1)D_y$ , given by  $g^y(W | x) = \text{diag}[\tilde{x} \tilde{x}^\top, \dots, \tilde{x} \tilde{x}^\top]$ , where  $\text{diag}(\cdot)$  means the (block) diagonal matrix constructed by the given matrix entries.

A *nonlinear* layer increments a linear layer by adding an element-wise activation function applied on  $W^\top \tilde{x}$ , and then randomized wrt the choice of the neuron. By definition 1, its RFIM is given by

$$g^y(W | x) = \text{diag}[\nu_f(w_1, x) \tilde{x} \tilde{x}^\top, \dots, \nu_f(w_m, x) \tilde{x} \tilde{x}^\top], \quad (2)$$

where  $\nu_f(w_i, x)$  is given in Subsec. 3.1.

A *softmax* layer, which often appears as the last layer of a MLP, is given by  $y \in \{1, \dots, m\}$ , where  $p(y) = \eta_y = \frac{\exp(w_y \tilde{x})}{\sum_{i=1}^m \exp(w_i \tilde{x})}$ . Its RFIM is a dense matrix given by

$$g^y(W) = \begin{bmatrix} (\eta_1 - \eta_1^2) \tilde{x} \tilde{x}^\top & \cdots & -\eta_1 \eta_m \tilde{x} \tilde{x}^\top \\ -\eta_2 \eta_1 \tilde{x} \tilde{x}^\top & \cdots & -\eta_2 \eta_m \tilde{x} \tilde{x}^\top \\ \vdots & \ddots & \vdots \\ -\eta_m \eta_1 \tilde{x} \tilde{x}^\top & \cdots & (\eta_m - \eta_m^2) \tilde{x} \tilde{x}^\top \end{bmatrix}.$$

Notice that its  $i$ 'th diagonal block  $(\eta_i - \eta_i^2) \tilde{x} \tilde{x}^\top$  resembles the RFIM of a single *sigm* neuron.



### 3.3. RFIM of Two Layers

By eq. (2), the one-layer RFIM is a product metric (Jost, 2011) and does not consider the *inter-neuron correlations*, which must be obtained by looking at a larger subsystem. Consider a two-layer model with stochastic output  $\mathbf{y}$  around the mean vector  $f(\mathbf{C}^\top \tilde{\mathbf{h}})$ , where  $\mathbf{h} = f(\mathbf{W}^\top \tilde{\mathbf{x}})$ . For simplicity, we ignore inter-layer correlations between the first layer and the second layer and focus on the inter-neuron correlations within the first layer. To do this, both  $\mathbf{x}$  and  $\mathbf{C}$  are considered as references to compute the RFIM of  $\mathbf{W}$ . By definition 1,  $g^{\mathbf{y}}(\mathbf{W} | \mathbf{x}, \mathbf{C}) = [\mathbf{G}_{ij}]_{D_h \times D_h}$  and each block has the form

$$\mathbf{G}_{ij} = \sum_{l=1}^{D_y} c_{il} c_{jl} \nu_f(\mathbf{c}_l, \mathbf{h}) \nu_f(\mathbf{w}_i, \mathbf{x}) \nu_f(\mathbf{w}_j, \mathbf{x}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top.$$

Now that we have the one-layer and two-layer RFIMs, we can either split a given feed-forward neural network into one-layer subsystems or into two-layer subsystems. A trade-off is that using a larger subsystem entails greater analytical and computational difficulty, although it could more accurately model the global system dynamics. In the extreme case, the FIM is obtained if the whole system is considered as one single subsystem.

## 4. RFIM: Key Advantages

This section discusses the theoretical advantages of the RFIM over the FIM. Consider wlog a MLP with Bernoulli outputs  $\mathbf{y} \in \{0, 1\}^m$ , whose mean  $\boldsymbol{\mu}$  is a deterministic function depending on the input  $\mathbf{x}$  and the network parameters  $\Theta$ . By Sec. 2, the FIM of the MLP can be computed as (see supplementary for proof)

$$\mathcal{I}(\Theta) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \frac{1}{\mu_j(\mathbf{x}_i)(1 - \mu_j(\mathbf{x}_i))} \frac{\partial \mu_j(\mathbf{x}_i)}{\partial \Theta} \frac{\partial \mu_j(\mathbf{x}_i)}{\partial \Theta^\top}. \quad (3)$$

Therefore  $\text{rank}(\mathcal{I}(\Theta)) \leq nm$ . The rank of a diagonal block of  $\mathcal{I}(\Theta)$  corresponding to one layer is even smaller. In a deep neural network (e.g. Szegedy, Christian et al. 2015), if the sample size  $n < \dim(\Theta)/m$ , then  $\mathcal{I}(\Theta)$  is doomed to be singular. All methods trying to approximate the FIM suffer from this problem and therefore rely on proper regularizations. If the network is decomposed into layers, the RFIM of each subsystem (layer) is given by eq. (2). Each sample can contribute maximally 1 to the rank of the neuron-RFIM and can contribute maximally  $D_y$  to the rank of the layer-RFIM. It only requires  $\max_i \{\dim(\mathbf{w}_i)\}$  (the maximum layer width) observations to have a full rank RFIM, where  $\mathbf{w}_i$  is the weight vector of the  $i$ 'th neuron. The RFIM is expected to have a much higher rank than the FIM. Higher rank means less singularity and more information is captured. Models that can

be distinguished by the RFIM may be identical in the sense of the FIM. Essentially, the RFIM integrates the *internal randomness* (Bengio, 2013) of the neural system by considering the output of each layer as a random variable. In theory, the FIM should also consider stochastic neurons. However it requires marginalizing the joint distribution of  $\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{y}$ . This makes the already infeasible computation even more challenging.

The RFIM is not an approximation of the FIM but is an *accurate* metric, defining the geometry of  $\theta$  wrt to its direct response  $\mathbf{h}$  in the system, or adjacent nodes in a graphical model. By the example in fig. 1,  $g^{\mathbf{y}}(\theta_L)$  of the last layer is exactly the corresponding block in  $\mathcal{I}(\Theta)$ : they both characterize how  $\theta_L$  affects the mapping  $\mathbf{h}_{L-1} \rightarrow \mathbf{y}$ . They start to diverge from the second to last layer. To compute the geometry of  $\theta_{L-1}$ , the RFIM looks at how  $\theta_{L-1}$  affects the local mapping  $\mathbf{h}_{L-2} \rightarrow \mathbf{h}_{L-1}$ , which can be measured reliably regardless of the rest of the system (think of a “debugging” process to separate and measure a single component). In contrast, the FIM examines how  $\theta_{L-1}$  affects the non-local mapping  $\mathbf{h}_{L-2} \rightarrow \mathbf{y}$ . This is a difficult task because it must consider the correlation between different layers. As an approximation, the block diagonalized version of the FIM ignores such correlations and therefore faces the loss of accuracy.

The RFIM makes it possible to maintain global system stability so that the intrinsic variations of different subsystems are balanced during learning. Consider a set of interconnected subsystems with internal parameters  $\{\theta_l\}$  and the corresponding response variables  $\{\mathbf{h}_l\}$ . The RFIM  $g^{\mathbf{h}_l}(\theta_l)$  measures how much the likelihood surface of  $\mathbf{h}_l$  is curved wrt a small learning step  $\delta\theta_l$ . By constraining the squared Riemannian distance  $\delta\theta_l^\top g^{\mathbf{h}_l}(\theta_l) \delta\theta_l$  having similar scales, different subsystems will present similar variations during learning. Within one subsystem, the learning along sensitive parameter directions is penalized. Among different subsystems, the learning of sensitive subsystems is penalized. Globally, the inter-subsystem stochastic connections have similar variance, maintaining a stable reference system and achieving efficient learning. This is similar to the idea of batch normalization (BN) (Ioffe & Szegedy, 2015) but has a deeper theoretical foundation.

Formally, we have the following theorem.

**Theorem 3.** Consider a learning system represented by a joint distribution  $p(\mathbf{x}, \mathbf{h})$  of  $\mathbf{x}$  (observables) and  $\mathbf{h}$  (hidden variables which connect subsystems). The joint FIM  $\mathcal{J}(\Theta) = E_p \left( \frac{\log p(\mathbf{x}, \mathbf{h} | \Theta)}{\partial \Theta} \frac{\log p(\mathbf{x}, \mathbf{h} | \Theta)}{\partial \Theta^\top} \right)$  has a block diagonal form. Each block is  $E_p(g^{\mathbf{h}}(\theta))$ , where  $\theta$  is the parameters within a subsystem and  $\mathbf{h}$  is its response variables to neighbour subsystems.

The global correspondence of the local RFIM is the joint

FIM. By theorem 3, the square distance  $d\Theta^\top \mathcal{J}(\Theta) d\Theta = E_p(\sum_l d\theta_l^\top g^{h_l}(\theta_l) d\theta_l)$  measures the system variance, including both the observables  $\mathbf{x}$  and the hidden variables  $\mathbf{h}$ . An intrinsic trade-off between the RFIM and the FIM is learning system stability versus efficiency. Normalizing the FIM is more efficient because it helps to achieve Fisher efficiency (Amari, 1998). Normalizing the RFIM is more stable since the hidden variations are bounded, which only guarantees subsystem Fisher efficiency characterized by the Cramér-Rao lower bound of local parameters.

## 5. Relative Natural Gradient Descent

The traditional *non-parametric* way of applying natural gradient requires re-calculating the FIM and solving a large linear system in each learning step. Besides the huge computational cost, it has a large approximation error. For example during online learning, a mini-batch of samples cannot faithfully reflect the “true” geometry, which has to integrate the risk of sample variations. That is, the FIM of a mini-batch is likely to be singular or poorly conditioned.

A recent series of efforts (Montavon & Müller, 2012; Raiko et al., 2012; Desjardins et al., 2015) are gearing towards a parametric approach to applying natural gradient, which memorizes and learns a geometry. For example, natural neural networks (Desjardins et al., 2015) augment each layer with a redundant linear layer, and let these linear layers parametrize the geometry of the neural manifold.

By dividing the learning system into subsystems, the RFIM potentially gives a systematical implementation of parametric natural gradient descent. The memory complexity of storing the Riemannian metric has been reduced from  $O(D^2)$  to  $O(\sum_i D_i^2)$ , where  $D_i = \dim(\mathbf{w}_i)$  is the size of the  $i$ 'th neuron. Consider there are  $M$  neurons in total, then the memory cost is reduced by a factor of  $M$ . The computational complexity has been reduced from  $O(D^e)$  ( $e \approx 2.373$ , Williams 2012) to  $O(\sum_i D_i^e)$ . Optimization based on RFIM is called Relative Natural Gradient Descent (RNGD).

The good performance of batch normalization (Ioffe & Szegedy, 2015) provides an empirical support for the RFIM. Basically, BN uses an *inter-sample* normalization layer to transform the layer input  $\mathbf{x}$  to  $\mathbf{z}$  with zero mean and unit variance and thus reduces “internal covariate shift”. In a typical case, above this normalization layer is a linear layer given by  $\mathbf{y} = \mathbf{W}^\top \mathbf{z}$ . If each dimension of  $\mathbf{z}$  is normalized, then the diagonal blocks of the linear layer RFIM  $g^{\mathbf{y}}(\mathbf{W}) = \text{diag}[\tilde{\mathbf{z}}\tilde{\mathbf{z}}^\top, \dots, \tilde{\mathbf{z}}\tilde{\mathbf{z}}^\top]$  become a covariance matrix with identity diagonal entries (after taking an empirical average). This gives the coordinate system  $\mathbf{W}$  a well conditioned RFIM for efficient learning.

### 5.1. RNGD with a relu MLP

This subsection builds a proof-of-concept experiment on MLP optimization. We partition the MLP into layers (one layer consists of a linear layer plus an element-wise non-linear activation function) as the subsystems. By eq. (2), the RFIM of layer  $l$  ( $l = 1, \dots, L$ ) with input  $\mathbf{h}_{l-1}$  ( $\mathbf{h}_0 = \mathbf{x}$ ) and weights  $\{\mathbf{w}_{l1}, \dots, \mathbf{w}_{lm_l}\}$  is

$$\text{diag} \left[ \nu_f(\mathbf{w}_{l1}, \mathbf{h}_{l-1}) \tilde{\mathbf{h}}_{l-1} \tilde{\mathbf{h}}_{l-1}^\top, \dots, \nu_f(\mathbf{w}_{lm_l}, \mathbf{h}_{l-1}) \tilde{\mathbf{h}}_{l-1} \tilde{\mathbf{h}}_{l-1}^\top \right].$$

The subsystem stability during one learning step  $\delta \mathbf{w}$  can be measured geometrically by  $\sum_{l=1}^L \sum_{i=1}^{m_l} \nu_f(\mathbf{w}_{li}, \mathbf{h}_{l-1}) (\delta \mathbf{w}_{li}^\top \tilde{\mathbf{h}}_{l-1})^2$ . Using this term as the geometric cost (the Lagrange term) in the trust region approach in Sec. 2, we get the following RNGD method. In a stochastic gradient descent scenario, each neuron  $i$  in layer  $l$  is updated by

$$\mathbf{w}_{li}^{\text{new}} \leftarrow \mathbf{w}_{li}^{\text{old}} - \mathbf{G}_{li}^{-1} \frac{\partial E}{\partial \mathbf{w}_{li}},$$

where  $E$  is the cost function and  $\mathbf{G}_{li}$  is a *learned metric*. The consideration is that a mini-batch of samples do not contain enough information to compute the RFIM, which should be averaged over all training samples. Therefore, for the  $i$ 'th neuron in layer  $l$ ,  $\mathbf{G}_{li}$  is initialized to identity, and is updated based on

$$\mathbf{G}_{li}^{\text{new}} \leftarrow (1 - \lambda) \mathbf{G}_{li}^{\text{old}} + \lambda \overline{\nu_f(\mathbf{w}_{li}, \mathbf{h}_{l-1}) \tilde{\mathbf{h}}_{l-1} \tilde{\mathbf{h}}_{l-1}^\top} + \epsilon \mathbf{I},$$

where  $\epsilon > 0$  is a hyper-parameter to avoid singularity caused by small sample size, and the average is taken over all samples in a mini-batch, and  $\lambda$  is a learning rate. In theory,  $\lambda$  should be gradually reduced to zero to guarantee the convergence of this geometry learning. To avoid solving a linear system in each iteration, every  $T$  iterations we recompute and store  $\mathbf{G}_{li}^{-1}$  based on the most updated  $\mathbf{G}_{li}$ . In the next  $T$  iterations, this  $\mathbf{G}_{li}^{-1}$  will be used as an approximation of the inverse RFIM. For the input layer which scales with the number of input features, and the final soft-max layer, we apply instead the RFIM of the corresponding linear layer to improve the computational efficiency.

We compare different optimizers on classifying MNIST digits. The network has shape 784-80-80-80-10, with relu activation units, a final soft-max layer, and uses the per-sample average cross-entropy with  $L_2$ -regularization as the learning cost function. We experiment on two different architectures: one is a plain MLP (PLAIN); the other has a batch normalization layer after each hidden layer (BNA), where a rescaling parameter is applied to ensure enough flexibility of the parametric structure (Ioffe & Szegedy, 2015). For simplicity, the architecture, mini-batch size (50), and  $L_2$  regularization strength ( $10^{-3}$ ) are fixed to be the same for all compared methods. The observations are consistent when these configurations vary.

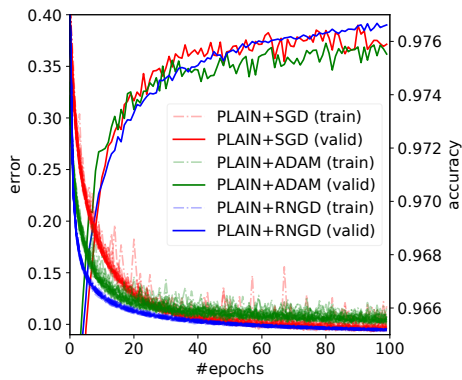
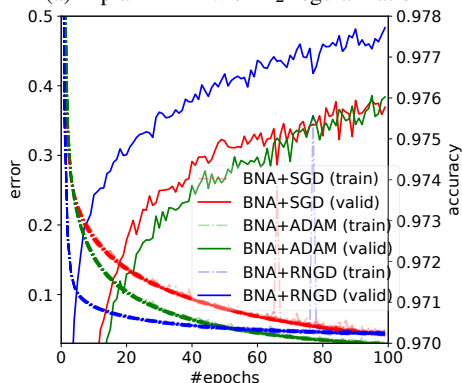

 (a) A plain MLP with  $L_2$  regularization

 (b) A MLP with batch normalization and  $L_2$  regularization

Figure 2. Learning curves of different optimizers on a MLP with two different architectures (with and without BN). The best learning rate for each method is selected based on the validation accuracy. Using this learning rate, the learning curves wrt 40 different random initializations are shown. The mean validation curve is shown for a clear visualization.

Figure 2 shows the learning curves of different methods. SGD is stochastic gradient descent. ADAM is the Adam optimizer (Kingma & Ba, 2014) with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . Our RNGD is implemented by modifying TensorFlow’s (Abadi, Martín et al., 2015) SGD optimizer. We set empirically  $T = 100$ ,  $\lambda = 0.005$  and  $\omega = 1$ .

RNGD presents a sharper learning curve and better generalization, especially when it is combined with BN. In this case, the final training error of RNGD is slightly larger than ADAM because by validation it favors a larger learning rate, which is applied on the neural network weights (based on RNGD) and BN parameters (based on SGD). For the ReLU activation,  $\nu_f(\mathbf{w}_i, \mathbf{x})$  is approximately binary, emphasizing such *informative samples* with  $\mathbf{w}_i^T \mathbf{x} > 0$ , which are the ones contributing to the learning of  $\mathbf{w}_i$  with non-zero gradient values. Each output neuron has a different subset of informative samples. RNGD normalizes  $\mathbf{x}$  differently wrt different output neurons, so that the in-

formative samples for each output neuron are centered and decorrelated.

In the above experiment, RNGD’s computational time per each epoch is roughly 4 ~ 10 times more than SGD and ADAM on a modern graphic card. Therefore in terms of wall clock time RNGD does not show advantages. This can be improved by more efficient implementations with low rank approximation techniques and early stopping. Our RNGD prototype hints at a promising direction to develop scalable 2<sup>nd</sup>-order deep learning optimizers based on the RFIM.

## 6. Related Works on FIM Diagonalization

One may ponder whether we can always find a suitable parameterization that yields a diagonal FIM that is straightforward to invert. This fundamental problem of parameter orthogonalization was first investigated by Jeffreys (1998) for decorrelating the *parameters of interest* from the *nuisance parameters*. Fisher diagonalization yields parameter orthogonalization (Cox & Reid, 1987), and is proved useful when estimating  $\hat{\Theta}$  using a maximum likelihood estimator (MLE) that is asymptotically normally distributed,  $\hat{\Theta}_n \sim G(\Theta, \mathcal{I}^{-1}(\Theta)/n)$ , and efficient since the variance of the estimator matches the Cramér-Rao lower bound. Using the chain rule, this amounts to find a suitable parameterization  $\Omega = \Omega(\Theta)$  satisfying

$$\sum_{i,j} E \left[ \frac{\partial^2 l}{\partial \Theta_i \partial \Theta_j} \right] \frac{\partial \Theta_i}{\partial \Omega_k} \frac{\partial \Theta_j}{\partial \Omega_l} = 0, \quad \forall k \neq l.$$

Thus in general, we end up with  $\binom{D}{2} = \frac{D(D-1)}{2}$  (non-linear) partial differential equations to satisfy (Huzurbazar, 1950). Therefore, in general there is no solution when  $\binom{D}{2} > D$ , that is when  $D > 3$ . When  $D = 2$ , the single differential equation is usually solvable and tractable, and the solution may not be unique: For example, Huzurbazar (1950) reports two orthogonalization schemes for the location-scale families  $\{\frac{1}{\sigma} p_0(\frac{x-\mu}{\sigma})\}$  that include the Gaussian family and the Cauchy family. Sometimes, the structure of the differential equation system yields a solution: For example, Jeffreys (1998) reported a parameter orthogonalization for Pearson’s distributions of type I which is of order  $D = 4$ . Cox and Reid (1987) further investigated this topic with application to conditional inference, and provide examples (including the Weibull distribution).

From the viewpoint of geometry, the FIM induces a Riemannian manifold with metric tensor  $g(\Theta) = \mathcal{I}(\Theta)$ . When the FIM may be degenerate, this yields a pseudo-Riemannian manifold (Thomas, 2014). In differential geometry, orthogonalization amounts to transforming the square length infinitesimal element  $g_{ij} d\Theta_i d\Theta_j$  of a Riemannian geometry into an orthogonal system  $\omega$  with match-

ing square length infinitesimal element  $\Omega_{ii}d\Omega_i^2$ . However, such a global orthogonal metric does not exist (Huzurbazar, 1950) when  $D > 3$  for an arbitrary metric tensor, although interesting Riemannian parameterization structures may be derived in Riemannian 4D geometry (Grant & Vickers, 2009).

For NEFs, the FIM can be made *block-diagonal* easily by using the *mixed coordinate system* (Amari, 2016)  $(\Theta_{1:k}, \mathbf{H}_{k+1:D})$ , where  $\mathbf{H} = E_p[t(\mathbf{x})] = \nabla F(\Theta)$  is the moment parameter, for any  $k \in \{1, \dots, D-1\}$ , where  $v_{b:e}$  denotes the subvector  $(v_b, \dots, v_e)^\top$  of  $v$ . The geometry of NEFs is a dually flat structure (Amari, 2016) induced by the convex mgf, the potential function. It defines a dual affine coordinate systems  $e^i = \partial_i = \frac{\partial}{\partial H_i}$  and  $e_j = \partial^j = \frac{\partial}{\partial \Theta^j}$  that are orthogonal:  $\langle e^i, e_j \rangle = \delta_j^i$ , where  $\delta_j^i = 1$  iff  $i = j$  and  $\delta_j^i = 0$  otherwise. Hence the FIM has two diagonal blocks. Those dual affine coordinate systems are defined up to an affine invertible transformation:  $\tilde{\Theta} = \mathbf{A}\Theta + \mathbf{b}$ ,  $\tilde{\mathbf{H}} = \mathbf{A}^{-1}\mathbf{H} + \mathbf{c}$ . In particular, for any order-2 NEF ( $D = 2$ ), we can *always* obtain two mixed parameterizations  $(\Theta_1, H_2)$  or  $(H_1, \Theta_2)$ .

The RFIM contributes another line of thought in parameter diagonalization. We investigate the Fisher information of *hidden variables*, or internal interfaces in the learning machine. This is novel since the majority of previous works concentrate on the FIM of the observables, or the external interface of the machine. From a causality perspective, we factor out the main cause (parameters within the subsystem) of the response variable with a direct action-reaction relationship, and regard the remaining parameters as a *reference* that can be easily estimated by the empirical distribution. This simplification may lead to broader applications of Fisher information in machine learning.

The particular case of a mixed coordinate system (that is not an affine coordinate system) induces in information geometry (Amari, 2016) a dual pair of orthogonal  $e$ - and  $m$ -orthogonal foliations. Our splits in RFIMs consider general non-orthogonal foliations that provide the factorization decompositions of the whole manifold into submanifolds, that are the leaves of the foliation (see section 3.7 of Amari & Nagaoka 2000).

## 7. Conclusion and Discussions

We investigate local structures of large learning systems using the new concept of Relative Fisher Information Metric. The key advantage of this approach is that the local learning dynamics can be analyzed in an accurate way without approximation. We present a core list of such local structures in neural networks, and give their corresponding RFIMs. This list of recipes can be used to provide guiding principles to design new optimizers for deep learning.

Our work applies to mirror descent as well since natural gradient is related to mirror descent (Raskutti & Mukherjee, 2015) as follows: In mirror descent to minimize a cost function  $E(\Theta)$ , given a strictly convex distance function  $D(\cdot, \cdot)$  in the first argument (playing the role of the proximity function), we express the gradient descent step as:

$$\Theta_{t+1} = \arg \min_{\Theta} \left\{ \Theta^\top \nabla E(\Theta_t) + \frac{1}{\gamma} D(\Theta, \Theta_t) \right\}.$$

When  $D(\Theta, \Theta')$  is chosen as a Bregman divergence  $B_F(\Theta, \Theta') = F(\Theta) - F(\Theta') - (\Theta - \Theta')^\top \nabla F(\Theta')$  wrt to a convex function  $F$ , it has been proved that the mirror descent on the  $\Theta$ -parameterization is equivalent (Raskutti & Mukherjee, 2015) to the natural gradient optimization on the induced Riemannian manifold with metric tensor  $(\nabla^2 F(\Theta))$  parameterized by the dual coordinate system  $\mathbf{H} = \nabla F(\Theta)$ .

In general, to perform a Riemannian gradient descent for minimizing a real-valued function  $f(\Theta)$  on the manifold, one needs to choose a proper metric tensor given in matrix form  $\mathbf{G}(\Theta)$ . Thomas (2014) constructed a toy example showing that the natural gradient may diverge while the ordinary gradient (for  $\mathbf{G} = \mathbf{I}$ ) converges. Recently, Thomas et al. (2016) proposed a new kind of descent method based on what they called the Energetic Natural Gradient that generalizes the natural gradient. The energy distance  $D_E(p(\Theta_1), p(\Theta_2))^2 = E[2d_{p(\Theta_1)}(X, Y) - d_{p(\Theta_1)}(X, X') - d_{p(\Theta_1)}(Y, Y')]$  where  $X, X' \sim p(\Theta_1)$  and  $Y, Y' \sim p(\Theta_2)$ , where  $d_{p(\Theta_1)}(\cdot, \cdot)$  is a distance metric over the support. Using a Taylor's expansion on their energy distance, they get the Energy Information Matrix (in a way similar to recovering the FIM from a Taylor's expansion of any  $f$ -divergence like the Kullback-Leibler divergence). Their idea is to incorporate prior knowledge on the structure of the support (observation space) to define energy distance. Twisting the geometry of the support (say, Wasserstein's optimal transport) with the geometry of the parametric distributions (Fisher-Rao geodesic distances) is indeed important (Chizat et al., 2015). In information geometry, invariance on the support is provided by a Markov morphism that is a probabilistic mapping of the support to itself (Čencov, 1982). There is no neighbourhood structure on the support in IG. Markov morphism includes deterministic transformation of a random variable by a statistic. It is well-known that  $\mathcal{I}_T(\Theta) \preceq \mathcal{I}_X(\Theta)$  with equality iff.  $T = T(X)$  is a sufficient statistic of  $X$ . Thus to get the same invariance for the energy distance (Thomas et al., 2016), one shall further require  $d_{p(\Theta)}(T(X), T(Y)) = d_{p(\Theta)}(X, Y)$ .

We believe that RFIMs will provide a sound methodology to build further efficient systems for deep learning. The full source codes to reproduce the experimental results are available at <https://www.lix.polytechnique.fr/~nielsen/RFIM>.



## Acknowledgements

The authors would like to thank the anonymous reviewers and Yann Ollivier for the helpful comments. This work was mainly conducted when the first author was a postdoctoral researcher at École Polytechnique.

## References

- Abadi, Martín et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from [tensorflow.org](https://www.tensorflow.org).
- Amari, Shun'ichi. Information geometry of the EM and em algorithms for neural networks. *Neural Networks*, 8(9): 1379–1408, 1995.
- Amari, Shun'ichi. Neural learning in structured parameter spaces – natural Riemannian gradient. In *NIPS 9*, pp. 127–133. MIT Press, 1997.
- Amari, Shun'ichi. Natural gradient works efficiently in learning. *Neural Comput.*, 10(2):251–276, 1998.
- Amari, Shun'ichi. *Information Geometry and its Applications*, volume 194 of *Applied Mathematical Sciences*. Springer Japan, 2016.
- Amari, Shun'ichi and Nagaoka, Hiroshi. *Methods of Information Geometry*, volume 191 of *Translations of Mathematical Monographs*. AMS and OUP, 2000. (Published in Japanese in 1993).
- Bengio, Yoshua. Estimating or propagating gradients through stochastic neurons. *CoRR*, abs/1305.2982, 2013.
- Bonnabel, Silvére. Stochastic gradient descent on Riemannian manifolds. *IEEE Trans. Automat. Contr.*, 58(9): 2217–2229, 2013.
- Čencov, Nikolai Nikolaevich. *Statistical decision rules and optimal inference*, volume 53 of *Translations of Mathematical Monographs*. American Mathematical Society, 1982. Translation from the Russian (published in 1972) edited by Lev J. Leifman.
- Chizat, Lenaïc, Schmitzer, Bernhard, Peyré, Gabriel, and Vialard, François-Xavier. An Interpolating Distance between Optimal Transport and Fisher-Rao. *arXiv e-prints*, 2015. 1506.06430 [math.AP].
- Clevert, Djork-Arné, Unterthiner, Thomas, and Hochreiter, Sepp. Fast and accurate deep network learning by exponential linear units (ELUs). *CoRR*, abs/1511.07289, 2015.
- Cobb, Loren, Koppstein, Peter, and Chen, Neng Hsin. Estimation and moment recursion relations for multimodal distributions of the exponential family. *JASA*, 78(381): 124–130, 1983.
- Cox, D. R. and Reid, N. Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society. Series B (Methodological)*, 49(1):1–39, 1987.
- Cramér, Harald. *Mathematical Methods of Statistics*, volume 9 of *Princeton Mathematical Series*. Princeton University Press, 1946.
- Desjardins, Guillaume, Simonyan, Karen, Pascanu, Razvan, and Kavukcuoglu, Koray. Natural neural networks. In *NIPS 28*, pp. 2071–2079. Curran Associates, Inc., 2015.
- Dowty, James G. Chentsov's theorem for exponential families. *arXiv preprints*, 2017. 1701.08895 [math.ST].
- Efron, Bradley and Hinkley, David V. Assessing the accuracy of the maximum likelihood estimator: Observed versus expected Fisher information. *Biometrika*, 65(3): 457–487, 1978.
- Fréchet, Maurice. Sur l'extension de certaines évaluations statistiques au cas de petits échantillons. *Revue de l'Institut International de Statistique / Review of the International Statistical Institute*, 11(3/4):182–205, 1943.
- Grant, James DE and Vickers, JA. Block diagonalization of four-dimensional metrics. *Classical and Quantum Gravity*, 26(23):235014, 2009.
- He, Kaiming, Zhang, Xiangyu, Ren, Shaoqing, and Sun, Jian. Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. In *ICCV*, 2015.
- Hinton, Geoffrey E., Osindero, Simon, and Teh, Yee-Whye. A fast learning algorithm for deep belief nets. *Neural Comput.*, 18(7):1527–1554, 2006.
- Hotelling, Harold. Spaces of statistical parameters. *American Mathematical Society Meeting*, 1929. (unpublished. Presented orally by O. Ore during the meeting).
- Huzurbazar, Vasant Shankar. Probability distributions and orthogonal parameters. *Mathematical Proceedings of the Cambridge Philosophical Society*, 46(02):281–284, 1950.
- Ioffe, Sergey and Szegedy, Christian. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML; JMLR: W&CP 37*, pp. 448–456, 2015.
- Jeffreys, Harold. *Theory of Probability*. Oxford Classic Texts in the Physical Sciences. OUP, 3rd edition, 1998. First published in 1939.

- Jost, Jürgen. *Riemannian Geometry and Geometric Analysis*. Springer, 6th edition, 2011.
- Kingma, Diederik P. and Ba, Jimmy. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- Kingma, Diederik P and Welling, Max. Auto-encoding variational Bayes. In *ICLR*, 2014. arXiv:1312.6114 [stat.ML].
- Kurita, Takio. Iterative weighted least squares algorithms for neural networks classifiers. *New Generation Computing*, 12(4):375–394, 1994.
- Le Roux, Nicolas, Manzagol, Pierre-Antoine, and Bengio, Yoshua. Topmoumoute online natural gradient algorithm. In *NIPS 20*, pp. 849–856. Curran Associates, Inc., 2008.
- Lebanon, Guy. *Riemannian geometry and statistical machine learning*. PhD thesis, CMU, 2005.
- Looks, Moshe, Herreshoff, Marcello, Hutchins, DeLesley, and Norvig, Peter. Deep learning with dynamic computation graphs. In *ICLR*, 2017. arXiv:1702.02181 [cs.NE].
- Marceau-Caron, Gaétan and Ollivier, Yann. Practical Riemannian neural networks. *CoRR*, abs/1602.08007, 2016.
- Martens, James. Deep learning via Hessian-free optimization. In *ICML*, pp. 735–742, 2010.
- Martens, James. New perspectives on the natural gradient method. *CoRR*, abs/1412.1193, 2014.
- Martens, James and Grosse, Roger. Optimizing neural networks with Kronecker-factored approximate curvature. In *ICML; JMLR: W&CP 37*, pp. 2408–2417, 2015.
- Montanari, Andrea. Computational implications of reducing data to sufficient statistics. *Electron. J. Statist.*, 9(2): 2370–2390, 2015.
- Montavon, Grégoire and Müller, Klaus-Robert. Deep Boltzmann machines and the centering trick. In *Neural Networks: Tricks of the Trade*, pp. 621–637. Springer Berlin Heidelberg, 2nd edition, 2012.
- Nair, Vinod and Hinton, Geoffrey E. Rectified linear units improve restricted Boltzmann machines. In *ICML*, pp. 807–814, 2010.
- Nielsen, Frank. Cramér-Rao lower bound and information geometry. *CoRR*, abs/1301.3578, 2013.
- Ollivier, Yann. Riemannian metrics for neural networks. *CoRR*, abs/1303.0818, 2013.
- Pascanu, Razvan and Bengio, Yoshua. Revisiting natural gradient for deep networks. In *ICLR*, 2014. arXiv:1301.3584 [cs.LG].
- Raiko, Tapani, Valpola, Harri, and LeCun, Yann. Deep learning made easier by linear transformations in perceptrons. In *AISTATS; JMLR W&CP 22*, pp. 924–932, 2012.
- Rao, Callyampudi Radhakrishna. Information and accuracy attainable in the estimation of statistical parameters. *Bull. Cal. Math. Soc.*, 37(3):81–91, 1945.
- Raskutti, Garvesh and Mukherjee, Sayan. The information geometry of mirror descent. In *Geometric Science of Information (GSI)*, volume 9389 of *LNCS*, pp. 359–368. Springer, 2015.
- Sun, Ke and Marchand-Maillet, Stéphane. An information geometry of statistical manifold learning. In *ICML; JMLR W&CP 32(2)*, pp. 1–9, 2014.
- Szegedy, Christian et al. Going deeper with convolutions. In *CVPR*, 2015.
- Thomas, Philip. GeNGA: A generalization of natural gradient ascent with positive and negative convergence results. In *ICML; JMLR W&CP 32*, pp. 1575–1583, 2014.
- Thomas, Philip, da Silva, B. C., Dann, C., and Brunskill, E. Energetic natural gradient descent. In *ICML*, 2016.
- Wager, Stefan, Wang, Sida, and Liang, Percy S. Dropout training as adaptive regularization. In *NIPS 26*, pp. 351–359. Curran Associates, Inc., 2013.
- Watanabe, Sumio. *Algebraic Geometry and Statistical Learning Theory*, volume 25 of *Cambridge Monographs on Applied and Computational Mathematics*. CUP, 2009.
- Williams, Virginia Vassilevska. Multiplying matrices faster than Coppersmith-Winograd. In *Annual ACM Symposium on Theory of Computing*, STOC’12, pp. 887–898, 2012.
- Zegers, Pablo. Fisher information properties. *Entropy*, 17: 4918–4939, 2015.