

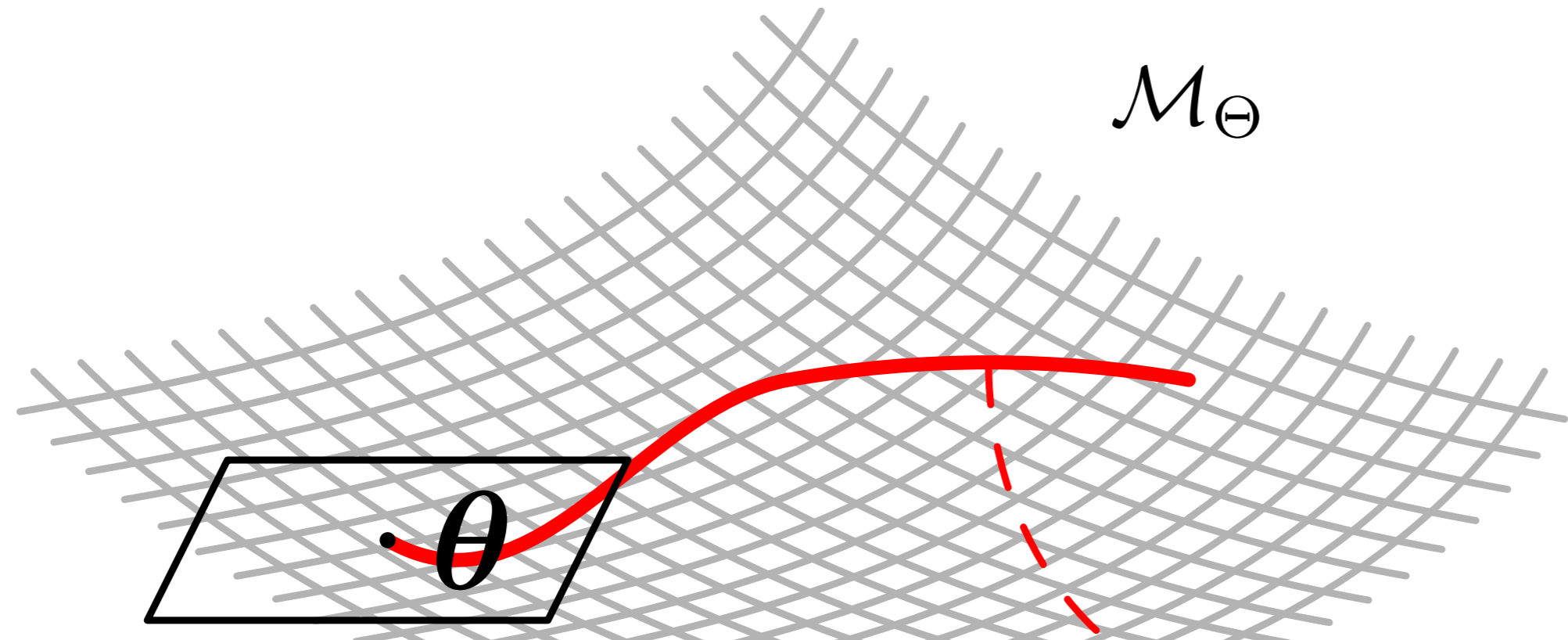
Fisher Information Metric (FIM)

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

$$\mathcal{I}_{ij} = E_p \left[\frac{\partial l}{\partial \Theta_i} \frac{\partial l}{\partial \Theta_j} \right] = -E_p \left[\frac{\partial^2 l}{\partial \Theta_i \partial \Theta_j} \right] = 4 \int \frac{\partial \sqrt{p(x | \Theta)}}{\partial \Theta_i} \frac{\partial \sqrt{p(x | \Theta)}}{\partial \Theta_j} dx, \quad (1)$$

where $l(\Theta) = \log p(x | \Theta)$ denotes the log-likelihood.

► Any parametric learning is inside a corresponding parameter manifold \mathcal{M}_Θ



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

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

FIM of a Multilayer Perceptron (MLP)

$$p(y | x, \Theta) = \sum_{h_1, \dots, h_{L-1}} p(y | h_{L-1}, \theta_L) \cdots p(h_2 | h_1, \theta_2) p(h_1 | x, \theta_1),$$

The FIM of a MLP has the following expression

$$g(\Theta) = E_{x \sim \hat{p}(X_n), y \sim p(y | x, \Theta)} \left[\frac{\partial l}{\partial \Theta} \frac{\partial l}{\partial \Theta^T} \right] = \frac{1}{n} \sum_{i=1}^n E_{p(y | x_i, \Theta)} \left[\frac{\partial l_i}{\partial \Theta} \frac{\partial l_i}{\partial \Theta^T} \right],$$

where $\hat{p}(X_n)$ is the empirical distribution of the samples $X_n = \{x_i\}_{i=1}^n$, and $l_i(\Theta) = \log p(y | x_i, \Theta)$ is the conditional log-likelihood.

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

$$\langle \delta\Theta, \delta\Theta \rangle_{g(\Theta)} = \delta\Theta^T g(\Theta) \delta\Theta = \frac{1}{n} \sum_{i=1}^n E_{p(y | x_i, \Theta)} \left[\delta\Theta^T \frac{\partial l_i}{\partial \Theta} \right]^2$$

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

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

$$\underbrace{L(\Theta_t + \delta\Theta)}_{\text{Loss function}} + \frac{1}{2\gamma} \underbrace{\langle \delta\Theta, \delta\Theta \rangle_{g(\Theta_t)}}_{\text{Squared step size}} \approx L(\Theta_t) + \delta\Theta^T \nabla L(\Theta_t) + \frac{1}{2\gamma} \delta\Theta^T g(\Theta_t) \delta\Theta,$$

giving a learning step

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

► Equivalence with mirror descent (Raskutti & Mukherjee 2013)

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 (RFIM) — Informal Ideas

- Decompose the learning system into subsystems
- The subsystems are interfaced with each other through hidden variables h_i
- Some subsystems are interfaced with the I/O environment through x_i and y_i
- Compute the subsystem FIM by **integrating out its interface variables h_i** , so that the intrinsics of this subsystem can be discussed regardless of the remaining parts



Given θ_f , how sensitive is r wrt tiny movements of θ ?

RFIM – Formal Definition

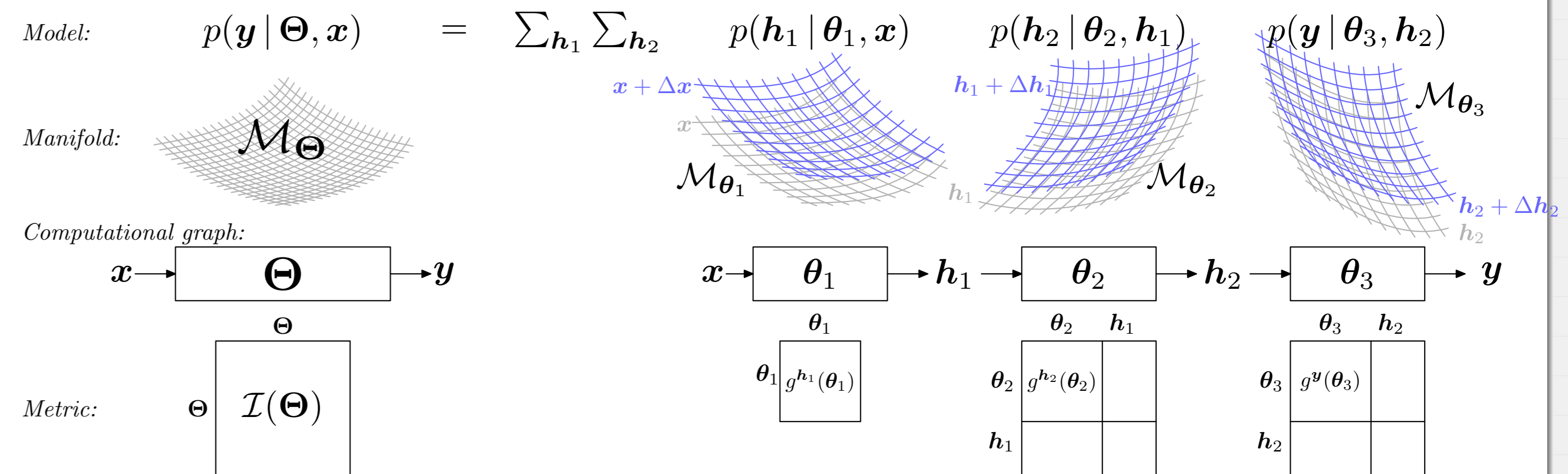
Given θ_f (the **reference**), the Relative Fisher Information Metric (RFIM) of θ wrt h (the **response**) is

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

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

- RFIM includes FIM as a special case.
- RFIM is dynamic wrt the reference θ_f

A Dynamic Geometry



► As the interface hidden variables 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 x , weights w , a hyperbolic tangent activation function, and a stochastic output $y \in \{-1, 1\}$, given by

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

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

$$g^y(w | x) = \nu_{\tanh}(w, x) \tilde{x} \tilde{x}^T, \quad \nu_{\tanh}(w, x) = \text{sech}^2(w^T \tilde{x}) = 1 - \tanh^2(w^T \tilde{x}).$$

Meaning: The RFIM has a large magnitude on the “learning zone” of the neuron.

A List of RFIMs

Subsystem the RFIM $g^y(w)$

- A tanh neuron $\text{sech}^2(w^T \tilde{x}) \tilde{x} \tilde{x}^T$
- A sigm neuron $\text{sigm}(w^T \tilde{x}) [1 - \text{sigm}(w^T \tilde{x})] \tilde{x} \tilde{x}^T$
- A relu neuron $[\iota + (1 - \iota) \text{sigm}(\frac{1-\iota}{\omega} w^T \tilde{x})]^2 \tilde{x} \tilde{x}^T$
- A elu neuron $\begin{cases} \tilde{x} \tilde{x}^T & \text{if } w^T \tilde{x} \geq 0 \\ (\alpha \exp(w^T \tilde{x}))^2 \tilde{x} \tilde{x}^T & \text{if } w^T \tilde{x} < 0 \end{cases}$
- A linear layer $\text{diag}[\tilde{x} \tilde{x}^T, \dots, \tilde{x} \tilde{x}^T]$
- A non-linear layer $\text{diag}[\nu_f(w_1, \tilde{x}) \tilde{x} \tilde{x}^T, \dots, \nu_f(w_m, \tilde{x}) \tilde{x} \tilde{x}^T]$
- A soft-max layer $\begin{bmatrix} (\eta_1 - \eta_1^2) \tilde{x} \tilde{x}^T & -\eta_1 \eta_2 \tilde{x} \tilde{x}^T & \dots & -\eta_1 \eta_m \tilde{x} \tilde{x}^T \\ -\eta_2 \eta_1 \tilde{x} \tilde{x}^T & (\eta_2 - \eta_2^2) \tilde{x} \tilde{x}^T & \dots & -\eta_2 \eta_m \tilde{x} \tilde{x}^T \\ \vdots & \vdots & \ddots & \vdots \\ -\eta_m \eta_1 \tilde{x} \tilde{x}^T & -\eta_m \eta_2 \tilde{x} \tilde{x}^T & \dots & (\eta_m - \eta_m^2) \tilde{x} \tilde{x}^T \end{bmatrix}$

Two layers see the paper.

Relative Natural Gradient Descent (RNGD)

For each subsystem,

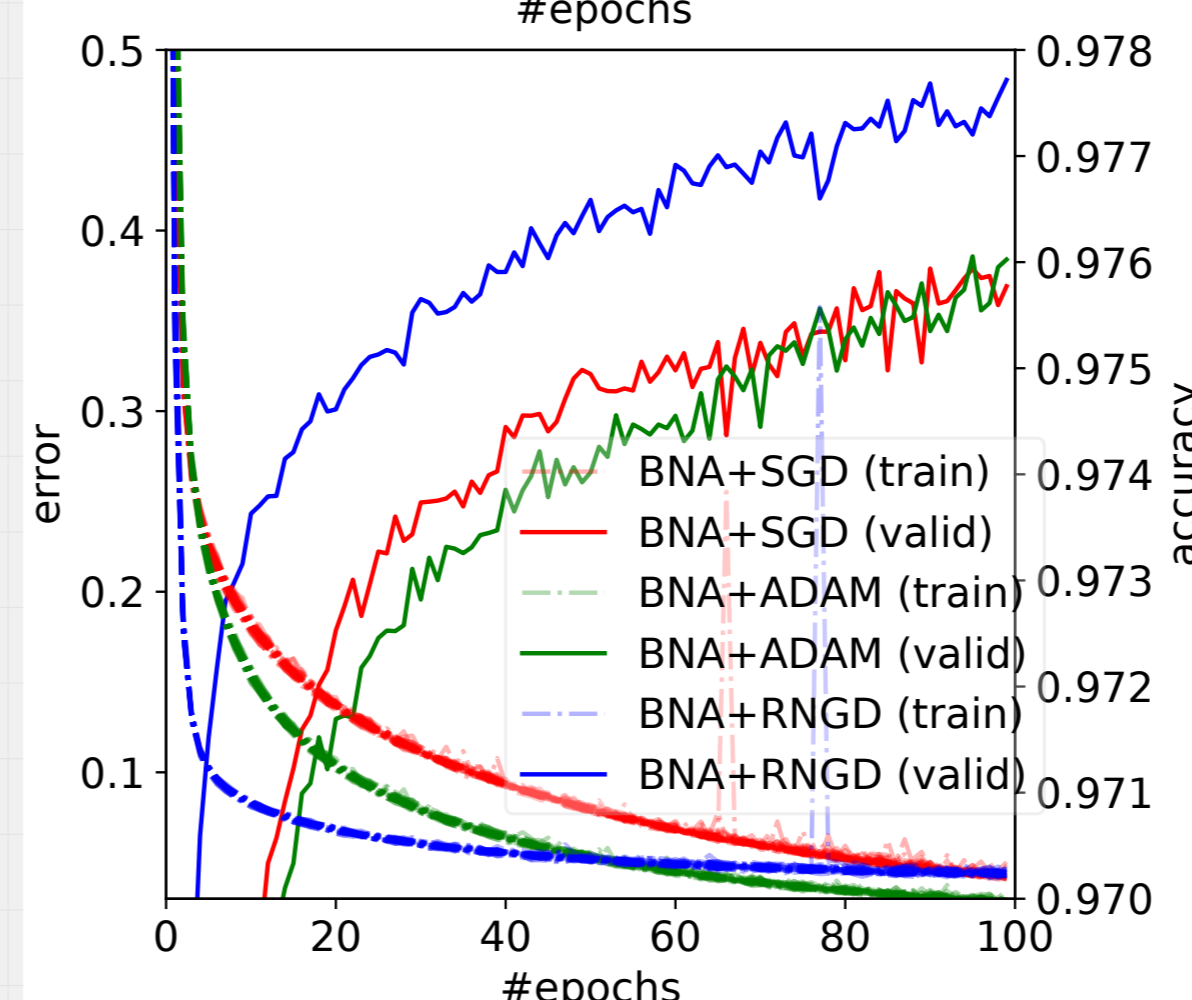
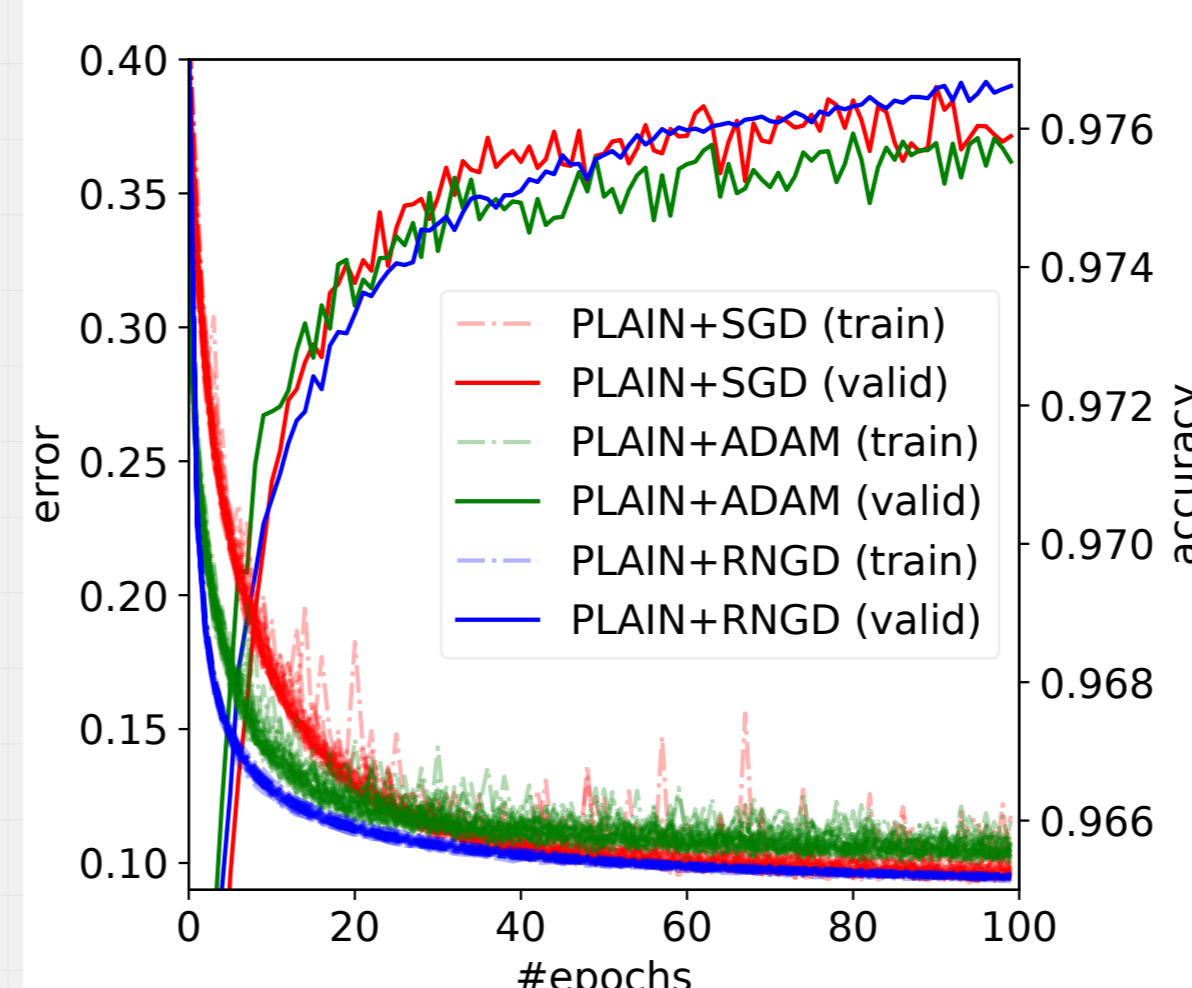
$$\theta_{t+1} \leftarrow \theta_t - \gamma \cdot \underbrace{(\bar{g}^h(\theta_t | \theta_f))^{-1}}_{\text{inverse RFIM}} \cdot \frac{\partial L}{\partial \theta} \Big|_{\theta=\theta_t}$$

where

$$\bar{g}^h(\theta_t | \theta_f) = \frac{1}{n} \sum_{i=1}^n g^h(\theta_t | \theta_f^i).$$

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

A Proof-of-concept



- MLP with shape 784-80-80-80-10
- relu activation
- Mini-batch size 50
- No dropout
- L_2 regularization
- Maintain an exponential moving average of the RFIM
- Recompute the inverse RFIM every 100 mini-batches
- PLAIN: a plain MLP
- BNA: a MLP with batch normalization
- Observations
- RNGD achieved sharper learning curve in terms of # iterations
- The computation cost of each epoch is several times more expensive
- RNGD can give better local optima

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
- Unlike NGD, RNGD can be implemented without approximation
- RFIM provides an accurate terminology to support feature whitening, natural neural networks, etc.