

Supplementary Material for “Relative Fisher Information and Natural Gradient for Learning Large Modular Models”

Ke Sun and Frank Nielsen

Contents

1	Non-linear Activation Functions	1
2	Examples of RFIMs	2
2.1	A Single <code>tanh</code> Neuron	3
2.2	A Single <code>sigm</code> Neuron	4
2.3	A Single <code>relu</code> Neuron	5
2.4	A Single <code>elu</code> Neuron	5
2.5	RFIM of a Linear Layer	5
2.6	RFIM of a Non-Linear Layer	6
2.7	RFIM of a Softmax Layer	7
2.8	RFIM of Two layers	7
3	Proof of Theorem 3	8
4	Experimental Settings & Zoomed Learning Curves	8

1 Non-linear Activation Functions

By definition,

$$\mathbf{tanh}(t) \stackrel{\text{def}}{=} \frac{\exp(t) - \exp(-t)}{\exp(t) + \exp(-t)}, \quad (1)$$

and

$$\mathbf{sech}(t) \stackrel{\text{def}}{=} \frac{2}{\exp(t) + \exp(-t)}.$$

It is easy to verify that

$$\mathbf{sech}^2(t) = [1 + \mathbf{tanh}(t)][1 - \mathbf{tanh}(t)] = 1 - \mathbf{tanh}^2(t).$$

By eq. (1),

$$\begin{aligned}\tanh'(t) &= \frac{\exp(t) + \exp(-t)}{\exp(t) + \exp(-t)} - \frac{\exp(t) - \exp(-t)}{[\exp(t) + \exp(-t)]^2} [\exp(t) - \exp(-t)] \\ &= \frac{[\exp(t) + \exp(-t)]^2 - [\exp(t) - \exp(-t)]^2}{[\exp(t) + \exp(-t)]^2} = \frac{4}{[\exp(t) + \exp(-t)]^2} = \operatorname{sech}^2(t).\end{aligned}$$

By definition,

$$\operatorname{sigm}(t) \stackrel{\text{def}}{=} \frac{1}{1 + \exp(-t)}.$$

Therefore

$$\operatorname{sigm}'(t) = -\frac{1}{[1 + \exp(-t)]^2} (-\exp(-t)) = \frac{\exp(-t)}{[1 + \exp(-t)]^2} = \operatorname{sigm}(t)[1 - \operatorname{sigm}(t)].$$

A smoothed version of the `relu` function is given by

$$\operatorname{relu}_\omega(t) \stackrel{\text{def}}{=} \omega \ln \left(\exp \left(\frac{\iota t}{\omega} \right) + \exp \left(\frac{t}{\omega} \right) \right),$$

where $\omega > 0$ and $0 \leq \iota < 1$. Then,

$$\begin{aligned}\operatorname{relu}'_\omega(t) &= \omega \frac{1}{\exp \left(\frac{\iota t}{\omega} \right) + \exp \left(\frac{t}{\omega} \right)} \left(\frac{\iota}{\omega} \exp \left(\frac{\iota t}{\omega} \right) + \frac{1}{\omega} \exp \left(\frac{t}{\omega} \right) \right) \\ &= \frac{\iota \exp \left(\frac{\iota t}{\omega} \right) + \exp \left(\frac{t}{\omega} \right)}{\exp \left(\frac{\iota t}{\omega} \right) + \exp \left(\frac{t}{\omega} \right)} \\ &= \iota + (1 - \iota) \frac{\exp \left(\frac{t}{\omega} \right)}{\exp \left(\frac{\iota t}{\omega} \right) + \exp \left(\frac{t}{\omega} \right)} \\ &= \iota + (1 - \iota) \frac{1}{\exp \left((\iota - 1) \frac{t}{\omega} \right) + 1} \\ &= \iota + (1 - \iota) \operatorname{sigm} \left(\frac{1 - \iota}{\omega} t \right).\end{aligned}\tag{2}$$

By definition,

$$\operatorname{elu}(t) = \begin{cases} t & \text{if } t \geq 0 \\ \alpha (\exp(t) - 1) & \text{if } t < 0. \end{cases}$$

Therefore

$$\operatorname{elu}'(t) = \begin{cases} 1 & \text{if } t \geq 0 \\ \alpha \exp(t) & \text{if } t < 0. \end{cases}\tag{3}$$

2 Examples of RFIMs

Table 1 shows a list of commonly used RFIMs, with detailed derivations given in the following subsections.

Table 1: Commonly used RFIMs

Subsystem	the RFIM $g^y(\mathbf{w})$
A tanh neuron	$\text{sech}^2(\mathbf{w}^\top \tilde{\mathbf{x}}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top$
A sigm neuron	$\text{sigm}(\mathbf{w}^\top \tilde{\mathbf{x}}) [1 - \text{sigm}(\mathbf{w}^\top \tilde{\mathbf{x}})] \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top$
A relu neuron	$[\iota + (1 - \iota) \text{sigm}(\frac{1-\iota}{\omega} \mathbf{w}^\top \tilde{\mathbf{x}})]^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top$
A elu neuron	$\begin{cases} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top & \text{if } \mathbf{w}^\top \tilde{\mathbf{x}} \geq 0 \\ (\alpha \exp(\mathbf{w}^\top \tilde{\mathbf{x}}))^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top & \text{if } \mathbf{w}^\top \tilde{\mathbf{x}} < 0 \end{cases}$
A linear layer	$\text{diag}[\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top, \dots, \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top]$
A non-linear layer	$\text{diag}[\nu_f(\mathbf{w}_1, \tilde{\mathbf{x}}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top, \dots, \nu_f(\mathbf{w}_m, \tilde{\mathbf{x}}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top]$
A soft-max layer	a dense matrix as shown in eq. (10)
Two layers	a dense matrix as shown in eq. (12)

2.1 A Single tanh Neuron

Consider a neuron with parameters \mathbf{w} and a Bernoulli output $y \in \{+, -\}$, $p(y = +) = p^+$, $p(y = -) = p^-$, and $p^+ + p^- = 1$. By the definition of RFIM, we have

$$\begin{aligned} g^y(\mathbf{w}) &= p^+ \frac{\partial \ln p^+}{\partial \mathbf{w}} \frac{\partial \ln p^+}{\partial \mathbf{w}^\top} + p^- \frac{\partial \ln p^-}{\partial \mathbf{w}} \frac{\partial \ln p^-}{\partial \mathbf{w}^\top} \\ &= \frac{1}{p^+} \frac{\partial p^+}{\partial \mathbf{w}} \frac{\partial p^+}{\partial \mathbf{w}^\top} + \frac{1}{p^-} \frac{\partial p^-}{\partial \mathbf{w}} \frac{\partial p^-}{\partial \mathbf{w}^\top}. \end{aligned}$$

Since $p^+ + p^- = 1$,

$$\frac{\partial p^+}{\partial \mathbf{w}} + \frac{\partial p^-}{\partial \mathbf{w}} = 0.$$

Therefore, the RFIM of a Bernoulli neuron has the general form

$$g^y(\mathbf{w}) = \left(\frac{1}{p^+} + \frac{1}{p^-} \right) \frac{\partial p^+}{\partial \mathbf{w}} \frac{\partial p^+}{\partial \mathbf{w}^\top} = \frac{1}{p^+ p^-} \frac{\partial p^+}{\partial \mathbf{w}} \frac{\partial p^+}{\partial \mathbf{w}^\top}. \quad (4)$$

A single **tanh** neuron with stochastic output $y \in \{-1, 1\}$ is given by

$$p(y = -1) = \frac{1 - \mu(\mathbf{x})}{2}, \quad (5)$$

$$p(y = 1) = \frac{1 + \mu(\mathbf{x})}{2}, \quad (6)$$

$$\mu(\mathbf{x}) = \text{tanh}(\mathbf{w}^\top \tilde{\mathbf{x}}). \quad (7)$$

By eq. (4),

$$\begin{aligned} g^y(\mathbf{w}) &= \frac{1}{\frac{1-\mu(\mathbf{x})}{2} \frac{1+\mu(\mathbf{x})}{2}} \left(\frac{1}{2} \frac{\partial \mu}{\partial \mathbf{w}} \right) \left(\frac{1}{2} \frac{\partial \mu}{\partial \mathbf{w}^\top} \right) \\ &= \frac{1}{(1 - \mu(\mathbf{x}))(1 + \mu(\mathbf{x}))} [1 - \mu^2(\mathbf{x})]^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \\ &= [1 - \mu^2(\mathbf{x})] \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \\ &= [1 - \text{tanh}^2(\mathbf{w}^\top \tilde{\mathbf{x}})] \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \\ &= \text{sech}^2(\mathbf{w}^\top \tilde{\mathbf{x}}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top. \end{aligned}$$

An alternative analysis is given as follows. By eqs. (5) to (7),

$$p(y = -1) = \frac{\exp(-\mathbf{w}^\top \tilde{\mathbf{x}})}{\exp(\mathbf{w}^\top \tilde{\mathbf{x}}) + \exp(-\mathbf{w}^\top \tilde{\mathbf{x}})},$$

$$p(y = 1) = \frac{\exp(\mathbf{w}^\top \tilde{\mathbf{x}})}{\exp(\mathbf{w}^\top \tilde{\mathbf{x}}) + \exp(-\mathbf{w}^\top \tilde{\mathbf{x}})}.$$

Then,

$$\begin{aligned} g^y(\mathbf{w}) &= E_{y \sim p(y|\mathbf{x})} \left(-\frac{\partial^2 \ln p(y)}{\partial \mathbf{w} \partial \mathbf{w}^\top} \right) \\ &= \frac{\partial^2}{\partial \mathbf{w} \partial \mathbf{w}^\top} \ln [\exp(\mathbf{w}^\top \tilde{\mathbf{x}}) + \exp(-\mathbf{w}^\top \tilde{\mathbf{x}})] \quad (\text{first linear term vanishes}) \\ &= \frac{\partial}{\partial \mathbf{w}^\top} \left[\frac{\exp(\mathbf{w}^\top \tilde{\mathbf{x}}) - \exp(-\mathbf{w}^\top \tilde{\mathbf{x}})}{\exp(\mathbf{w}^\top \tilde{\mathbf{x}}) + \exp(-\mathbf{w}^\top \tilde{\mathbf{x}})} \right] \tilde{\mathbf{x}} \\ &= \frac{\partial}{\partial \mathbf{w}^\top} \tanh(\mathbf{w}^\top \tilde{\mathbf{x}}) \tilde{\mathbf{x}} \\ &= \text{sech}^2(\mathbf{w}^\top \tilde{\mathbf{x}}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top. \end{aligned}$$

The intuitive meaning of $g^y(\mathbf{w})$ is a weighted covariance to emphasize such “informative” \mathbf{x} ’s that

- are in the linear region of `tanh`
- contain “ambiguous” samples

We will need at least $\dim(\mathbf{w})$ samples to make $g^y(\mathbf{w})$ full rank.

2.2 A Single `sigm` Neuron

A single `sigm` neuron is given by

$$\begin{aligned} p(y = 0) &= 1 - \mu(\mathbf{x}), \\ p(y = 1) &= \mu(\mathbf{x}), \\ \mu(\mathbf{x}) &= \text{sigm}(\mathbf{w}^\top \tilde{\mathbf{x}}). \end{aligned}$$

By eq. (4),

$$\begin{aligned} g^y(\mathbf{w}) &= \frac{1}{p(y = 0)p(y = 1)} \frac{\partial p(y = 1)}{\partial \mathbf{w}} \frac{\partial p(y = 1)}{\partial \mathbf{w}^\top} \\ &= \frac{1}{\mu(\mathbf{x})(1 - \mu(\mathbf{x}))} \frac{\partial \mu}{\partial \mathbf{w}} \frac{\partial \mu}{\partial \mathbf{w}^\top} \\ &= \frac{1}{\mu(\mathbf{x})(1 - \mu(\mathbf{x}))} \mu^2(\mathbf{x})(1 - \mu(\mathbf{x}))^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \\ &= \mu(\mathbf{x})(1 - \mu(\mathbf{x})) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \\ &= \text{sigm}(\mathbf{w}^\top \tilde{\mathbf{x}}) [1 - \text{sigm}(\mathbf{w}^\top \tilde{\mathbf{x}})] \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top. \end{aligned}$$

2.3 A Single relu Neuron

Consider a single neuron with Gaussian output $p(y | \mathbf{w}, \mathbf{x}) = G(y | \mu(\mathbf{w}, \mathbf{x}), \sigma^2)$. Then

$$\begin{aligned}
g^y(\mathbf{w} | \mathbf{x}) &= E_{p(y | \mathbf{w}, \mathbf{x})} \left[\frac{\partial \ln G(y | \mu, \sigma^2)}{\partial \mathbf{w}} \frac{\partial \ln G(y | \mu, \sigma^2)}{\partial \mathbf{w}^\top} \right] \\
&= E_{p(y | \mathbf{w}, \mathbf{x})} \left[\frac{\partial}{\partial \mathbf{w}} \left(-\frac{1}{2\sigma^2} (y - \mu)^2 \right) \frac{\partial}{\partial \mathbf{w}^\top} \left(-\frac{1}{2\sigma^2} (y - \mu)^2 \right) \right] \\
&= E_{p(y | \mathbf{w}, \mathbf{x})} \left[\left(-\frac{1}{\sigma^2} (\mu - y) \right)^2 \frac{\partial \mu}{\partial \mathbf{w}} \frac{\partial \mu}{\partial \mathbf{w}^\top} \right] \\
&= \frac{1}{\sigma^4} E_{p(y | \mathbf{w}, \mathbf{x})} (\mu - y)^2 \frac{\partial \mu}{\partial \mathbf{w}} \frac{\partial \mu}{\partial \mathbf{w}^\top} \\
&= \frac{1}{\sigma^2} \frac{\partial \mu}{\partial \mathbf{w}} \frac{\partial \mu}{\partial \mathbf{w}^\top}.
\end{aligned}$$

We set $\sigma = 1$ to get rid of a scale parameter of the RFIM. We get

$$g^y(\mathbf{w} | \mathbf{x}) = \frac{\partial \mu}{\partial \mathbf{w}} \frac{\partial \mu}{\partial \mathbf{w}^\top}. \quad (8)$$

A single **relu** neuron is given by

$$\mu(\mathbf{w}, \mathbf{x}) = \text{relu}_\omega(\mathbf{w}^\top \tilde{\mathbf{x}}).$$

By eqs. (2) and (8),

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

2.4 A Single elu Neuron

Similar to the analysis in Subsec. 2.3, a single **elu** neuron is given by

$$\mu(\mathbf{w}, \mathbf{x}) = \text{elu}(\mathbf{w}^\top \tilde{\mathbf{x}}).$$

By eq. (3),

$$\frac{\partial \mu}{\partial \mathbf{w}} = \begin{cases} \tilde{\mathbf{x}} & \text{if } \mathbf{w}^\top \tilde{\mathbf{x}} \geq 0 \\ \alpha \exp(\mathbf{w}^\top \tilde{\mathbf{x}}) \tilde{\mathbf{x}} & \text{if } \mathbf{w}^\top \tilde{\mathbf{x}} < 0. \end{cases}$$

By eq. (8),

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

2.5 RFIM of a Linear Layer

Consider a linear layer

$$p(\mathbf{y}) = G(\mathbf{y} | \mathbf{W}^\top \tilde{\mathbf{x}}, \sigma^2 \mathbf{I}),$$

where $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_{D_y})$. By the definition of the multivariate Gaussian distribution,

$$\ln p(\mathbf{y}) = -\frac{1}{2} \ln 2\pi - \frac{D_y}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^{D_y} (y_i - \mathbf{w}_i^\top \tilde{\mathbf{x}})^2.$$

Therefore,

$$\forall i, \quad \frac{\partial}{\partial \mathbf{w}_i} \ln p(\mathbf{y}) = -\frac{1}{\sigma^2} (\mathbf{w}_i^\top \tilde{\mathbf{x}} - y_i) \tilde{\mathbf{x}}.$$

Therefore,

$$\forall i, \forall j \quad \frac{\partial}{\partial \mathbf{w}_i} \ln p(\mathbf{y}) \frac{\partial}{\partial \mathbf{w}_j^\top} \ln p(\mathbf{y}) = \frac{1}{\sigma^4} (y_i - \mathbf{w}_i^\top \tilde{\mathbf{x}}) (y_j - \mathbf{w}_j^\top \tilde{\mathbf{x}}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top.$$

\mathbf{W} is vectorized by stacking its columns $\{\mathbf{w}_i\}_{i=1}^{D_y}$. In the following \mathbf{W} will be used interchangeably to denote either the matrix or its vector form. Correspondingly, the RFIM $g^{\mathbf{y}}(\mathbf{W})$ has $D_y \times D_y$ blocks, where the off-diagonal blocks are

$$\forall i \neq j, \quad E_{p(\mathbf{y})} \left(\frac{\partial}{\partial \mathbf{w}_i} \ln p(\mathbf{y}) \frac{\partial}{\partial \mathbf{w}_j^\top} \ln p(\mathbf{y}) \right) = \frac{1}{\sigma^4} E_{p(\mathbf{y})} [(y_i - \mathbf{w}_i^\top \tilde{\mathbf{x}}) (y_j - \mathbf{w}_j^\top \tilde{\mathbf{x}})] \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top = \mathbf{0},$$

and the diagonal blocks are

$$\forall i, \quad E_{p(\mathbf{y})} \left(\frac{\partial}{\partial \mathbf{w}_i} \ln p(\mathbf{y}) \frac{\partial}{\partial \mathbf{w}_i^\top} \ln p(\mathbf{y}) \right) = \frac{1}{\sigma^4} E_{p(\mathbf{y})} (y_i - \mathbf{w}_i^\top \tilde{\mathbf{x}})^2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top = \frac{1}{\sigma^2} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top.$$

In summary,

$$g^{\mathbf{y}}(\mathbf{W}) = \frac{1}{\sigma^2} \text{diag} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top, \dots, \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top].$$

By setting $\sigma = 1$ we get

$$g^{\mathbf{y}}(\mathbf{W}) = \text{diag} [\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top, \dots, \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top].$$

2.6 RFIM of a Non-Linear Layer

The statistical model of a non-linear layer with independent output units is

$$p(\mathbf{y} | \mathbf{W}, \mathbf{x}) = \prod_{i=1}^{D_y} p(y_i | \mathbf{w}_i, \mathbf{x}).$$

Then,

$$\ln p(\mathbf{y} | \mathbf{W}, \mathbf{x}) = \sum_{i=1}^{D_y} \ln p(y_i | \mathbf{w}_i, \mathbf{x}).$$

Therefore,

$$\frac{\partial^2}{\partial \mathbf{W} \partial \mathbf{W}^\top} \ln p(\mathbf{y} | \mathbf{W}, \mathbf{x}) = \begin{bmatrix} \frac{\partial^2}{\partial \mathbf{w}_1 \partial \mathbf{w}_1^\top} \ln p(y_1 | \mathbf{w}_1, \mathbf{x}) & & \\ & \ddots & \\ & & \frac{\partial^2}{\partial \mathbf{w}_{D_y} \partial \mathbf{w}_{D_y}^\top} \ln p(y_{D_y} | \mathbf{w}_{D_y}, \mathbf{x}) \end{bmatrix}.$$

Therefore the RFIM $g^{\mathbf{y}}(\mathbf{W})$ is a block-diagonal matrix, with the i 'th block given by

$$-E_{p(\mathbf{y} | \mathbf{W}, \mathbf{x})} \left[\frac{\partial^2}{\partial \mathbf{w}_i \partial \mathbf{w}_i^\top} \ln p(y_i | \mathbf{w}_i, \mathbf{x}) \right] = -E_{p(y_i | \mathbf{w}_i, \mathbf{x})} \left[\frac{\partial^2}{\partial \mathbf{w}_i \partial \mathbf{w}_i^\top} \ln p(y_i | \mathbf{w}_i, \mathbf{x}) \right],$$

which is simply the single neuron RFIM of the i 'th neuron.

2.7 RFIM of a Softmax Layer

Recall that

$$\forall i \in \{1, \dots, m\}, \quad p(y = i) = \frac{\exp(\mathbf{w}_i \tilde{\mathbf{x}})}{\sum_{i=1}^m \exp(\mathbf{w}_i \tilde{\mathbf{x}})}.$$

Then

$$\forall i, \quad \ln p(y = i) = \mathbf{w}_i \tilde{\mathbf{x}} - \ln \sum_{i=1}^m \exp(\mathbf{w}_i \tilde{\mathbf{x}}).$$

Hence

$$\forall i, \forall j, \quad \frac{\partial \ln p(y = i)}{\partial \mathbf{w}_j} = \delta_{ij} \tilde{\mathbf{x}} - \frac{\exp(\mathbf{w}_j \tilde{\mathbf{x}})}{\sum_{i=1}^m \exp(\mathbf{w}_i \tilde{\mathbf{x}})} \tilde{\mathbf{x}},$$

where $\delta_{ij} = 1$ if and only if $i = j$ and $\delta_{ij} = 0$ otherwise. Then

$$\begin{aligned} \forall i, \forall j, \forall k, \quad \frac{\partial^2 \ln p(y = i)}{\partial \mathbf{w}_j \partial \mathbf{w}_k^\top} &= -\delta_{jk} \frac{\exp(\mathbf{w}_j \tilde{\mathbf{x}})}{\sum_{i=1}^m \exp(\mathbf{w}_i \tilde{\mathbf{x}})} \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top + \frac{\exp(\mathbf{w}_j \tilde{\mathbf{x}})}{(\sum_{i=1}^m \exp(\mathbf{w}_i \tilde{\mathbf{x}}))^2} \exp(\mathbf{w}_k \tilde{\mathbf{x}}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \\ &= (-\delta_{jk} \eta_j + \eta_j \eta_k) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top. \end{aligned} \quad (9)$$

The right-hand-side of eq. (9) does not depend on i . Therefore

$$g^{\mathbf{y}}(\mathbf{W}) = \begin{bmatrix} (\eta_1 - \eta_1^2) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top & -\eta_1 \eta_2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top & \cdots & -\eta_1 \eta_m \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \\ -\eta_2 \eta_1 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top & (\eta_2 - \eta_2^2) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top & \cdots & -\eta_2 \eta_m \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \\ \vdots & \vdots & \ddots & \vdots \\ -\eta_m \eta_1 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top & -\eta_m \eta_2 \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top & \cdots & (\eta_m - \eta_m^2) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top \end{bmatrix}. \quad (10)$$

2.8 RFIM of Two layers

Consider a two layer structure, where the output \mathbf{y} satisfies a multivariate Bernoulli distribution with independent dimensions. By a similar analysis to Subsec. 2.1, we have

$$g^{\mathbf{y}}(\mathbf{W}) = \sum_{l=1}^{D_y} \nu_f(\mathbf{c}_l, \mathbf{h}) \frac{\partial \mathbf{c}_l^\top \mathbf{h}}{\partial \mathbf{W}} \frac{\partial \mathbf{c}_l^\top \mathbf{h}}{\partial \mathbf{W}^\top}. \quad (11)$$

It can be written block by block as $g^{\mathbf{y}}(\mathbf{W}) = [\mathbf{G}_{ij}]_{D_h \times D_h}$, where each block \mathbf{G}_{ij} means the correlation between the i 'th hidden neuron with weights \mathbf{w}_i and the j 'th hidden neuron with weights \mathbf{w}_j . By eq. (11),

$$\begin{aligned} \mathbf{G}_{ij} &= \sum_{l=1}^{D_y} \nu_f(\mathbf{c}_l, \mathbf{h}) \frac{\partial \mathbf{c}_l^\top \mathbf{h}}{\partial \mathbf{w}_i} \frac{\partial \mathbf{c}_l^\top \mathbf{h}}{\partial \mathbf{w}_j^\top} = \sum_{l=1}^{D_y} \nu_f(\mathbf{c}_l, \mathbf{h}) \frac{\partial c_{il} h_i}{\partial \mathbf{w}_i} \frac{\partial c_{jl} h_j}{\partial \mathbf{w}_j^\top} \\ &= \sum_{l=1}^{D_y} \nu_f(\mathbf{c}_l, \mathbf{h}) c_{il} c_{jl} \frac{\partial h_i}{\partial \mathbf{w}_i} \frac{\partial h_j}{\partial \mathbf{w}_j^\top} = \sum_{l=1}^{D_y} \nu_f(\mathbf{c}_l, \mathbf{h}) c_{il} c_{jl} (\nu_f(\mathbf{w}_i, \mathbf{x}) \tilde{\mathbf{x}}) (\nu_f(\mathbf{w}_j, \mathbf{x}) \tilde{\mathbf{x}}^\top) \\ &= \sum_{l=1}^{D_y} c_{il} c_{jl} \nu_f(\mathbf{c}_l, \mathbf{h}) \nu_f(\mathbf{w}_i, \mathbf{x}) \nu_f(\mathbf{w}_j, \mathbf{x}) \tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top. \end{aligned} \quad (12)$$

The proof of the other case, where two `relu` layers have stochastic output \mathbf{y} satisfying a multivariate Gaussian distribution with independent dimensions, is very similar and is omitted.

3 Proof of Theorem 3

Proof. By assumption, the joint distribution $p(\mathbf{x}, \mathbf{h})$ is in a factorable form. Therefore

$$\log p(\mathbf{x}, \mathbf{h}) = \sum_{l=1}^L \log p(\mathbf{h}_l | \boldsymbol{\theta}_l, \mathbf{r}_l), \quad (13)$$

where $l = 1, \dots, L$ is the index of subsystems, \mathbf{h}_l is the subsystem output, and \mathbf{r}_l is the reference of the subsystem. We have $\biguplus_{l=1}^L \{\mathbf{h}_l\} = \{\mathbf{x}, \mathbf{h}\}$ and $\biguplus_{l=1}^L \{\boldsymbol{\theta}_l\} = \{\boldsymbol{\Theta}\}$. Therefore

$$\begin{aligned} E_p \left(-\frac{\partial^2}{\partial \boldsymbol{\theta}_l \partial \boldsymbol{\theta}_l^\top} \log p(\mathbf{x}, \mathbf{h}) \right) &= E_p \left(-\frac{\partial^2}{\partial \boldsymbol{\theta}_l \partial \boldsymbol{\theta}_l^\top} \log p(\mathbf{h}_l | \boldsymbol{\theta}_l, \mathbf{r}_l) \right) \\ &= E_{p(\mathbf{r}_l)} \left(E_{p(\mathbf{h}_l | \mathbf{r}_l)} \left(-\frac{\partial^2}{\partial \boldsymbol{\theta}_l \partial \boldsymbol{\theta}_l^\top} \log p(\mathbf{h}_l | \boldsymbol{\theta}_l, \mathbf{r}_l) \right) \right) \\ &= E_p (g^{\mathbf{h}_l}(\boldsymbol{\theta}_l)), \end{aligned}$$

and

$$E_p \left(-\frac{\partial^2}{\partial \boldsymbol{\theta}_{l_1} \partial \boldsymbol{\theta}_{l_2}^\top} \log p(\mathbf{x}, \mathbf{h}) \right) = \mathbf{0} \quad (\forall l_1 \neq l_2).$$

Based on the Hessian expression of RFIM, $\mathcal{J}(\boldsymbol{\Theta})$ is in a block-diagonal form, with each block given by $E_p (g^{\mathbf{h}_l}(\boldsymbol{\theta}_l))$. \square

4 Experimental Settings & Zoomed Learning Curves

The training/validation/testing sets have 50,000/10,000/10,000 images, respectively. Each sample is a gray scale image of size 28×28 (784 dimensional feature space) and is labeled as one of ten different classes. For all methods, the mini-batch size is fixed to 50 and the L_2 regularization strength is fixed to 10^{-3} . For each optimizer, we try to find the best learning rate in the range $\{\dots, 10^{-1}, 5 \times 10^{-2}, 10^{-2}, 5 \times 10^{-3}, 10^{-3}, \dots\}$. On the tested architectures, a good learning rate configuration for RNGD is usually around 10^{-2} or 5×10^{-3} . The optimizers are in their default settings in TensorFlow 1.0. For the Adam optimizer, $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$. For RNGD, we set empirically $T = 100$, $\lambda = 0.005$ and $\omega = 1$. We use the Glorot uniform initializer to set the initial weights.

For each method and each learning rate configuration, we try 40 independent runs with different random seeds. Then, we select the best configuration based on the validation accuracy. Then, we plot the 40 learning curves as well as the average validation curve. The learning curves are obtained by evaluating the training error and validation accuracy after each epoch (one pass over all available training data).

See the following figs. (1–4) for the learning curves on four different architectures with **relu** activation units and L_2 regularization. Only the training curves and validation curves are shown for a clear presentation. The testing accuracy is close to the validation accuracy (run our codes to see the detailed results).

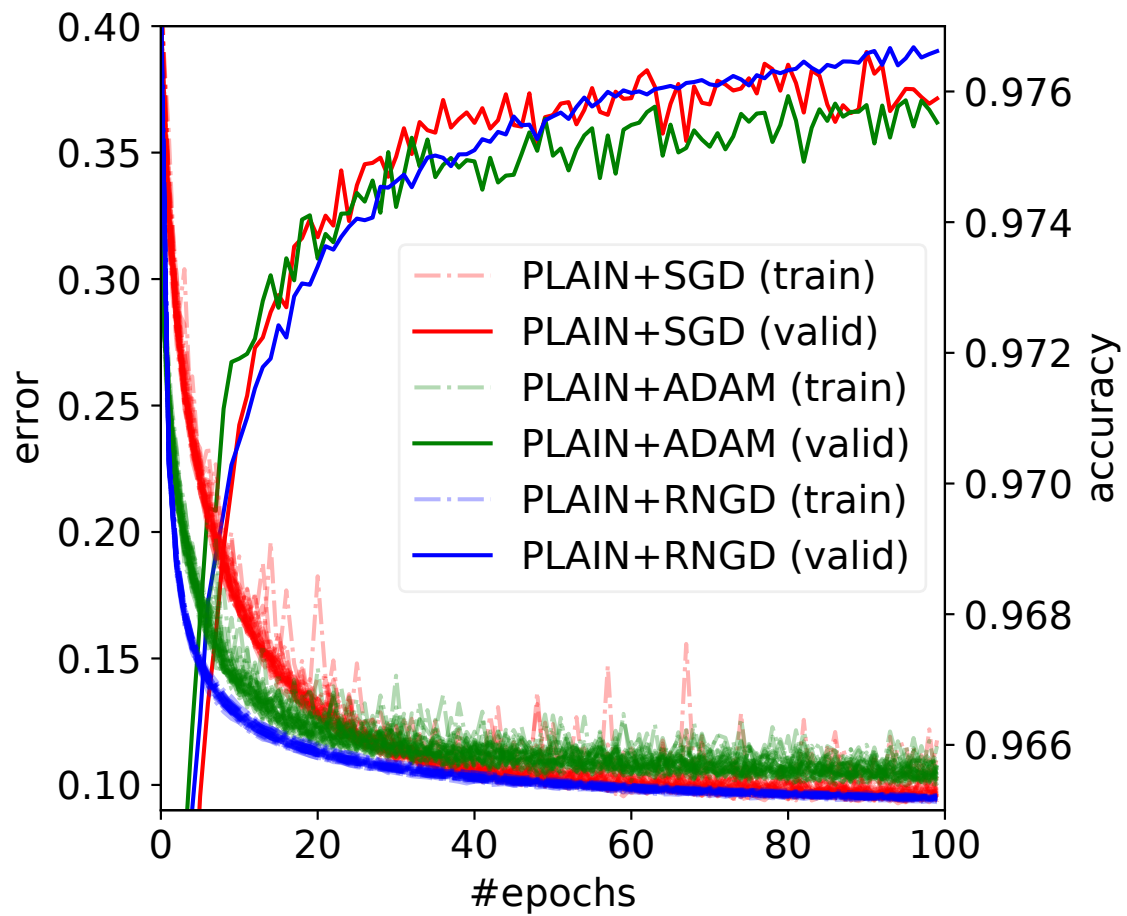


Figure 1: A MLP with shape 784-80-80-80-10.

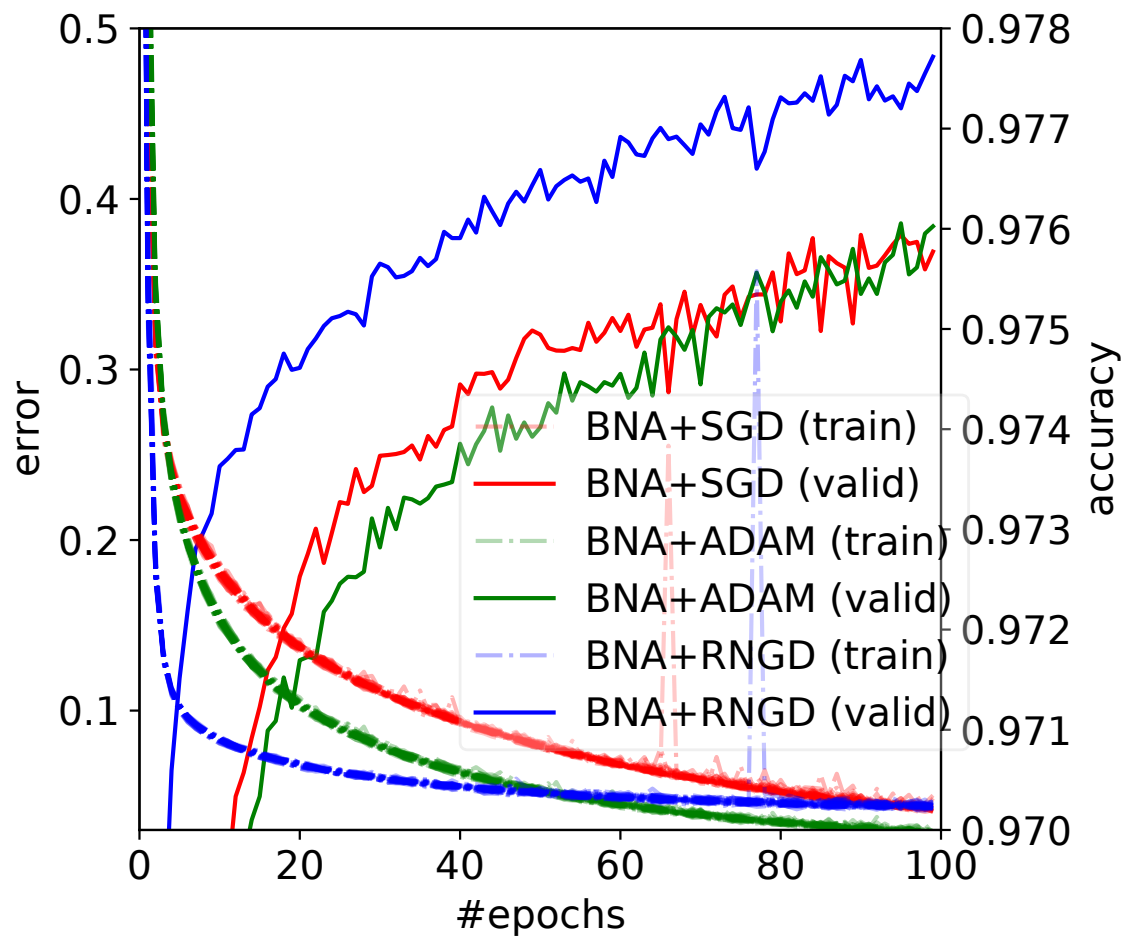


Figure 2: A MLP with shape 784–80–80–80–10 and batch normalization after each hidden layer.

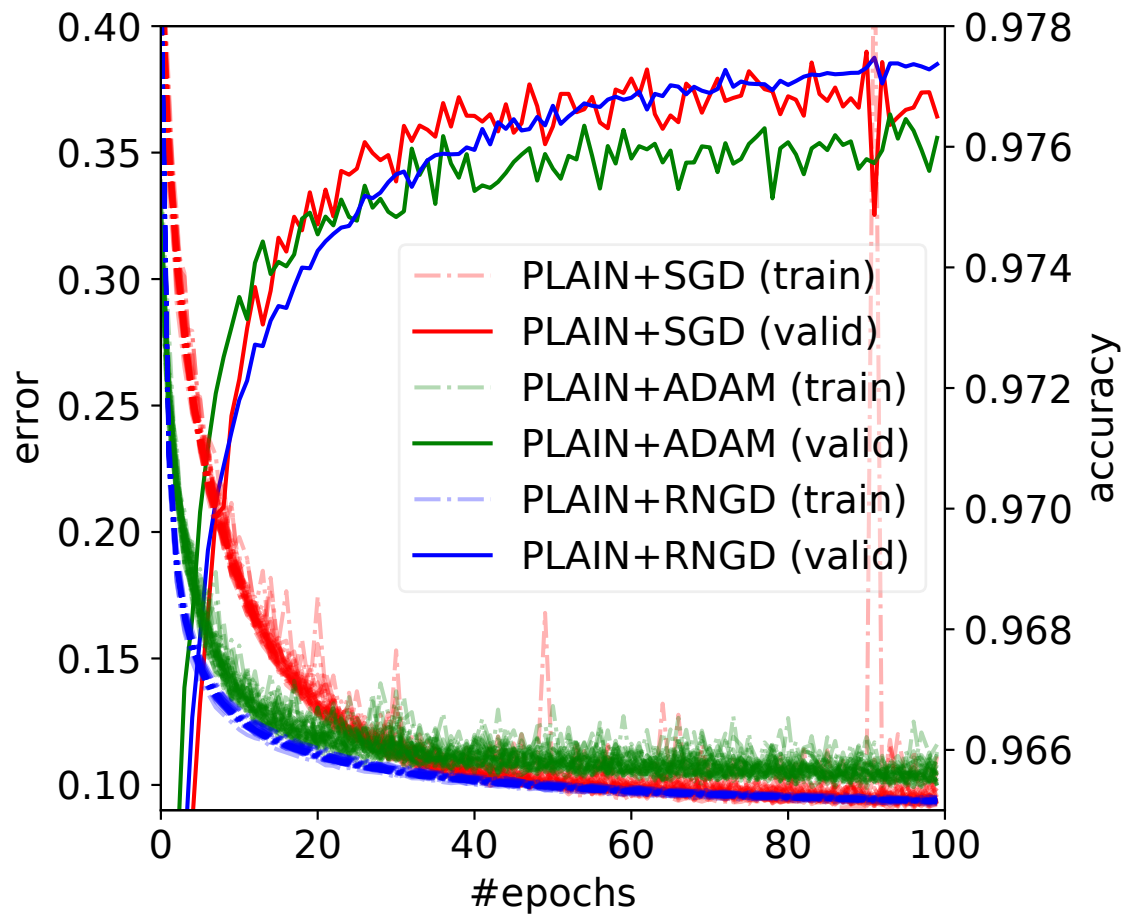


Figure 3: A MLP with shape 784-100-100-100-10.

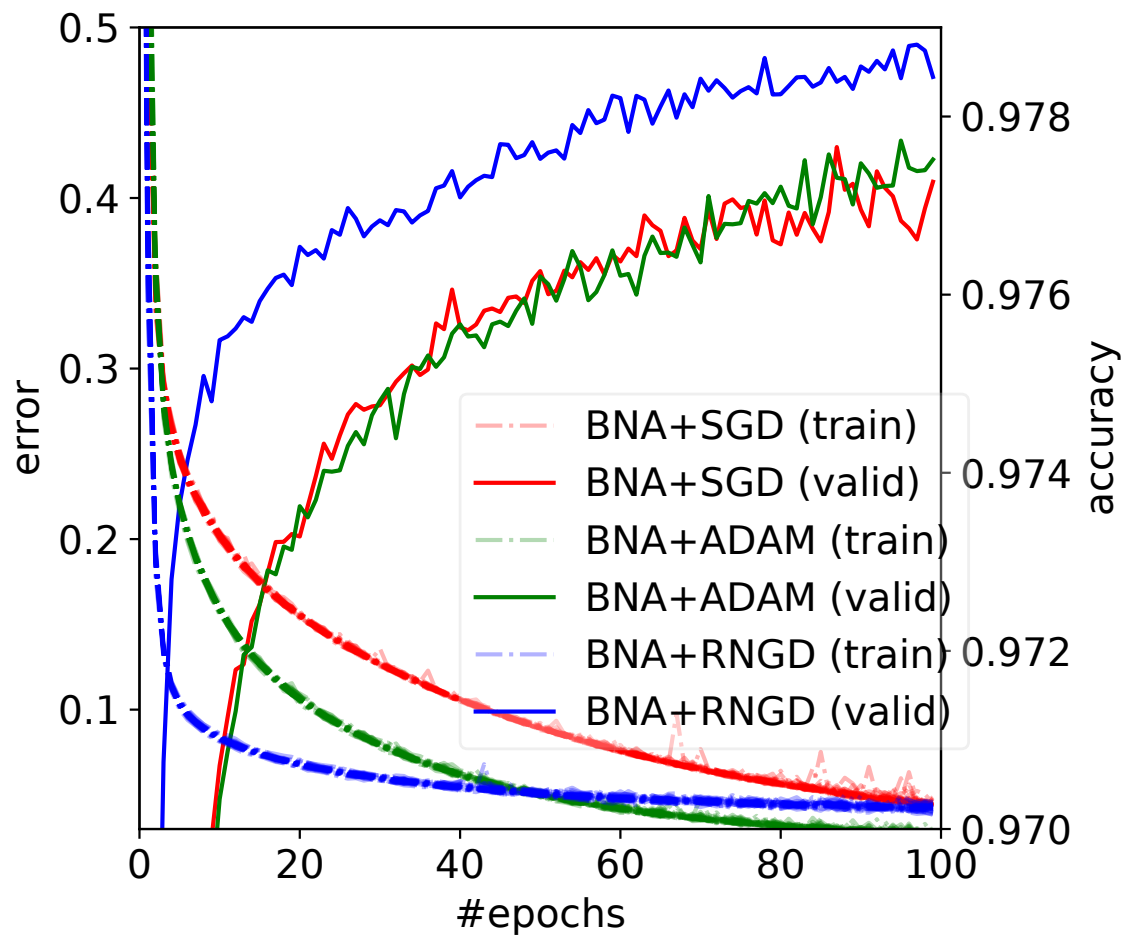


Figure 4: A MLP with shape 784–100–100–100–10 and batch normalization after each hidden layer.