

Uncertainty Quantification for Large-Scale Deep Neural Networks via Post-StoNet Modeling

University of Pennsylvania and Purdue University

Supplementary Material

S1 Theoretical Proofs

S1.1 Proof of Lemma 2

Proof. For simplicity of notations, we suppress the iteration index t . Let $\tilde{\mathbf{Y}}_l = \mathbf{b}_l + \mathbf{w}_l \Psi(\mathbf{Y}_{l-1})$ for $l = 2, \dots, h$, and let $\tilde{\mathbf{Y}}_1 = \mathbf{b}_1 + \mathbf{w}_1 \mathbf{X}$. By the definition of the StoNet model (2.2), \mathbf{Y}_l can be written as $\mathbf{Y}_l = \tilde{\mathbf{Y}}_l + \mathbf{e}_l$ for $l \in \{1, 2, \dots, h\}$.

Since σ_l^2 has been set to a very small value, we have $\Psi(\mathbf{Y}_l) \approx \Psi(\tilde{\mathbf{Y}}_l) + \nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l) \circ \mathbf{e}_l$, where \circ denotes elementwise product. Then

$$\begin{aligned} \Sigma_l &\approx \text{Var}(\mathbb{E}(\Psi(\tilde{\mathbf{Y}}_l) + \nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l) \circ \mathbf{e}_l | \tilde{\mathbf{Y}}_l)) + \mathbb{E}(\text{Var}(\Psi(\tilde{\mathbf{Y}}_l) + \nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l) \circ \mathbf{e}_l | \tilde{\mathbf{Y}}_l)) \\ &= \text{Var}(\Psi(\tilde{\mathbf{Y}}_l)) + \text{diag} \left\{ \sigma_l^2 \mathbb{E}[\nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l) \circ \nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l)] \right\}, \end{aligned} \tag{S1.1}$$

where $\text{diag}\{\mathbf{v}\}$ with $\mathbf{v} \in \mathbb{R}^d$ denotes a $d \times d$ diagonal matrix with diagonal elements being \mathbf{v} .

By Assumption 3-(iii), the activation function is bounded. For example, *tanh* or *sigmoid* is used in the model. By Assumption 1, there exists some constant C_1 such that $\|\mathbf{b}_l\|_\infty < C_1, \|\mathbf{w}_l\|_\infty < C_1$. By Assumption 3, $\|\mathbf{X}\|_\infty$ is bounded. Therefore, there exists some constant C_2 such that for any $L \in \{1, 2, \dots, h\}$, $\|\tilde{\mathbf{Y}}_l\|_\infty \leq C_1 + C_1 C_2$ holds by rescaling \mathbf{X} by a factor of $\prod_{l=1}^h d_l$. Since both $\Psi(\tilde{\mathbf{Y}}_l)$ and $\nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l)$ are bounded, there exists a constant $\kappa_{\max, l}$ such that

$$\phi_{\max}(d_{l,n}|\Sigma_l) \leq \kappa_{\max, l}.$$

To establish the lower bound, we note that $\|\tilde{\mathbf{Y}}_l\|_\infty \leq C_1 + C_1 C_2$. Therefore, for an activation function which has nonzero gradients on any closed interval, e.g., *tanh* and *sigmoid*, there exists a constant $C_3 > 0$ such that $\min_{i=1, \dots, d_l} \nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l)_i > C_3$, where $\nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l)_i$ denotes the i -th element of $\nabla_{\tilde{\mathbf{Y}}_l} \Psi(\tilde{\mathbf{Y}}_l)$. Then we can take $\kappa_{\min, l} = \sigma_{l,n}^2 C_3^2$ such that

$$\phi_{\min}(d_{l,n}|\Sigma_l) \geq \kappa_{\min, l},$$

which completes the proof. \square

S1.2 Proof of Part (i) of Theorem 1

Proof. By Lemma 2, $\Sigma_l^{(t)}$ satisfies the requirements of Theorem 1 of Meinshausen and Yu (2009) and Theorem 1 of Huang et al. (2008). Then, by Theorem 1 of Meinshausen and Yu (2009) (for linear regression) and Theorem 1 of Huang et al. (2008) (for logistic regression), we have r_n as given in the lemma by summarizing the l_2 -errors of coefficient estimation for all $\sum_{l=1}^{h+1} d_l$ regression/logistic regressions. Further, by the setting of $(\sigma_{1,n}^2, \dots, \sigma_{h+1,n}^2)$ as specified in Assumption 3, we have $r_n \rightarrow 0$ as $n \rightarrow \infty$. This completes the proof of part (i) of Theorem 1. \square

S1.3 Proof of Part (ii) of Theorem 1

Proof. Then part (ii) of Theorem 1 directly follows from Theorem 4 of Liang et al. (2018) that the estimator $\hat{\boldsymbol{\theta}}_n^{(t)}$ is consistent when both n and t are sufficiently large. \square

S1.4 Proof of Corollary 1

Proof. Let $\hat{\boldsymbol{\theta}}_n^{(t)}$ denote the estimate of $\boldsymbol{\theta}_n$ at iteration t , and let $\boldsymbol{\theta}_*^{(t)}$ denote its “true” value at iteration t , and let $\boldsymbol{\theta}^*$ denote its true value in the StoNet. By the proof of Theorem 4 of Liang et al. (2018) and Theorem 1 of Meinshausen and Yu (2009), for the StoNet with the linear regression output layer, we

have

$$\mathbb{E}\|\hat{\boldsymbol{\theta}}_n^{(t)} - \boldsymbol{\theta}^*\| \leq \frac{1}{1 - \rho^*} \mathbb{E}\|\hat{\boldsymbol{\theta}}_n^{(t)} - \boldsymbol{\theta}_*^{(t)}\| \prec \frac{\sqrt{r_n}}{1 - \rho^*}, \quad \text{as } t \rightarrow \infty, \quad (\text{S1.2})$$

by summarizing all $d_1 + d_2 + \cdots + d_{h+1}$ linear regressions, where ρ^* is a constant as defined in Assumption 5. For the StoNet with the logistic regression output layer, we have the same result by Theorem 1 of Huang et al. (2008). Further, by Markov inequality, there exists a constant c such that

$$P\left(\|\hat{\boldsymbol{\theta}}_n^{(t)} - \boldsymbol{\theta}^*\| > c\sqrt{r_n}\right) \rightarrow 0, \quad \text{as } n \rightarrow \infty \text{ and } t \rightarrow \infty.$$

Then, by Assumption 6,

- For any $i \in \gamma^*$, $\|\hat{\boldsymbol{\theta}}_n^{(t)} - \boldsymbol{\theta}^*\| \leq c\sqrt{r_n}$ implies $|\hat{\boldsymbol{\theta}}_{i,n}^{(t)}| > c\sqrt{r_n}$.
- For any $i \notin \gamma^*$, $\|\hat{\boldsymbol{\theta}}_n^{(t)} - \boldsymbol{\theta}^*\| \leq c\sqrt{r_n}$ implies $|\hat{\boldsymbol{\theta}}_{i,n}^{(t)}| < c\sqrt{r_n}$.

Therefore,

$$P(\hat{\gamma} = \gamma^*) \geq P(\|\hat{\boldsymbol{\theta}}_n^{(t)} - \boldsymbol{\theta}^*\| \leq c\sqrt{r_n}) \rightarrow 1, \quad \text{as } n \rightarrow \infty \text{ and } t \rightarrow \infty, \quad (\text{S1.3})$$

which concludes the proof. \square

S1.5 Proof of Corollary 2

Proof. This proof involves several notations, including $\widehat{\Sigma}$, $\hat{\theta}$, and $\varsigma_{h+1,j}$. As noted in the main text, their dependence on the sample size n is implicit and has been depressed for notational simplicity. As implied by (S4.9)-(S4.11), we have $\widehat{\Sigma}_i \rightarrow 0$, $i \in \{1, 2, \dots, h\}$, as $n \rightarrow \infty$. Additionally, as $n \rightarrow \infty$,

$$\|\mu(\mathbf{z}, \hat{\theta}) - \mu(\mathbf{z}, \theta^*)\| \xrightarrow{p} 0,$$

for any test point \mathbf{z} , and

$$\varsigma_{h+1,j}^2 - \sigma_{h+1}^2 \xrightarrow{p} 0, \quad j \in \{1, 2, \dots, d_{h+1}\}.$$

Therefore, the nominal coverage rate $1 - \alpha$ is asymptotically guaranteed as $n \rightarrow \infty$. □

S2 The Imputation Regularized-Optimization Algorithm for StoNet Training

This algorithm is given in Algorithm S1.

Algorithm S1: IRO Algorithm for StoNet

Input: Dataset $D_n = (\mathbb{Y}, \mathbb{X})$, total iteration number T , and Monte Carlo step number t_{MC} .

Initialization: Randomly initialize the network parameters

$$\hat{\boldsymbol{\theta}}^{(0)} = (\hat{\boldsymbol{\theta}}_1^{(0)}, \dots, \hat{\boldsymbol{\theta}}_{h+1}^{(0)}).$$

for $t = 1, 2, \dots, T$ **do**

• **Imputation:** For each sample $(\mathbf{X}^{(i)}, \mathbf{Y}^{(i)})$, draw $\mathbf{Y}_{\text{mis}}^{(i,t)}$ from $\pi(\mathbf{Y}_{\text{mis}}^{(i)} | \mathbf{Y}^{(i)}, \mathbf{X}^{(i)}, \hat{\boldsymbol{\theta}}_n^{(t-1)}, \boldsymbol{\sigma}_n^2)$ with a Metropolis or Langevin dynamics kernel by iterating for t_{MC} steps.

• **Regularized optimization:** Based on the pseudo-complete data $\{(\mathbf{Y}^{(i)}, \mathbf{Y}_{\text{mis}}^{(i,t)}, \mathbf{X}^{(i)}) : i = 1, 2, \dots, n\}$, update $\hat{\boldsymbol{\theta}}_n^{(t-1)}$ by minimizing a penalized loss function, i.e., setting

$$\hat{\boldsymbol{\theta}}_n^{(t)} = \arg \min_{\boldsymbol{\theta}} \left\{ -\frac{1}{n} \sum_{i=1}^n \log \pi(\mathbf{Y}^{(i)}, \mathbf{Y}_{\text{mis}}^{(i,t)} | \mathbf{X}^{(i)}, \boldsymbol{\theta}, \boldsymbol{\sigma}_n^2) + P_{\lambda_n}(\boldsymbol{\theta}) \right\}, \quad (\text{S2.4})$$

where the penalty $P_{\lambda_n}(\boldsymbol{\theta})$ is chosen such that $\hat{\boldsymbol{\theta}}_n^{(t)}$ forms a consistent estimator of

$$\begin{aligned} \boldsymbol{\theta}_*^{(t)} &= \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\boldsymbol{\theta}_n^{(t-1)}} \log \pi(\mathbf{Y}, \mathbf{Y}_{\text{mis}} | \mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}_n^2) \\ &= \arg \max_{\boldsymbol{\theta}} \int \log \pi(\mathbf{Y}_{\text{mis}}, \mathbf{Y} | \mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}_n^2) \pi(\mathbf{Y}_{\text{mis}} | \mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}_n^{(t-1)}, \boldsymbol{\sigma}_n^2) \\ &\quad \times \pi(\mathbf{Y} | \mathbf{X}, \boldsymbol{\theta}_*, \boldsymbol{\sigma}_n^2) d\mathbf{Y}_{\text{mis}} d\mathbf{Y}, \end{aligned} \quad (\text{S2.5})$$

where $\boldsymbol{\theta}_*^{(t)}$ is called the working true parameter at iteration t .

Output: $\hat{\boldsymbol{\theta}}_n^{(T)}$

S3 Adaptive Stochastic Gradient MCMC for Efficient StoNet Training

The IRO algorithm requires computation on the full dataset at each iteration and, therefore, it is less scalable with respect to big data. In practice, we can train the sparse StoNet using the adaptive stochastic gradient MCMC algorithm as proposed in (Liang et al., 2022). To make the paper self-contained, we give a review of the adaptive stochastic gradient MCMC algorithm below.

Let $\pi(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}^2) = \int \pi(\mathbf{Y}, \mathbf{Y}_{\text{mis}}|\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}^2)d\mathbf{Y}_{\text{mis}}$ denote the likelihood function of the observed data for the StoNet model. By Fisher's identity, we have

$$\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}^2) = \int \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{Y}, \mathbf{Y}_{\text{mis}}|\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}^2)\pi(\mathbf{Y}_{\text{mis}}|\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}^2)d\mathbf{Y}_{\text{mis}},$$

which implies the sparse StoNet can also be trained by solving the equation

$$\int \nabla_{\boldsymbol{\theta}} [\log \pi(\mathbf{Y}, \mathbf{Y}_{\text{mis}}|\mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}^2) + \log P_{\lambda}(\boldsymbol{\theta})]\pi(\mathbf{Y}_{\text{mis}}|\mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}, \boldsymbol{\sigma}^2)d\mathbf{Y}_{\text{mis}} = 0, \quad (\text{S3.6})$$

where $P_{\lambda}(\boldsymbol{\theta})$ denotes a penalty function satisfying Assumption 4. By Theorem 1 of Liang et al. (2018), solving (S3.6) will lead to the same solution

as solving the optimization problem specified below:

$$\hat{\boldsymbol{\theta}}_n^* = \arg \max_{\boldsymbol{\theta}} \left\{ \frac{1}{n} \sum_{i=1}^n \log \pi(\mathbf{Y}^{(i)}, \mathbf{Y}_{\text{mis}}^{(i)} | \mathbf{X}^{(i)}, \boldsymbol{\theta}, \sigma^2) + \frac{1}{n} P_{\lambda}(\boldsymbol{\theta}) \right\}.$$

By Deng et al. (2019) and Liang et al. (2022), the equation (S3.6) can be solved using an adaptive stochastic gradient MCMC algorithm, which works by iterating between the following two steps:

- (a) (*Sampling*) Generate $\mathbf{Y}_{\text{mis}}^{(k+1)}$ from a transition kernel induced by a stochastic gradient MCMC algorithm, e.g., stochastic gradient Hamilton Monte Carlo (SGHMC) (Chen et al., 2014).
- (b) (*Parameter updating*) Set $\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} + \gamma_{k+1} g(\mathbf{Y}_{\text{mis}}^{(k+1)}, U_{k+1})$, where γ_{k+1} denotes the step size used in the stochastic approximation procedure.

The pseudo-code of the adaptive SGHMC algorithm is given by Algorithm S2, where we let $\boldsymbol{\theta}_i = (\mathbf{w}_i, \mathbf{b}_i)$ denote the parameters associated with layer i of the StoNet, let $(\mathbf{Y}_0^{(s,k)}, \mathbf{Y}_{h+1}^{(s,k)}) = (\mathbf{X}^{(s)}, \mathbf{Y}^{(s)})$ denote a training sample s , and let $\mathbf{Y}_{\text{mis}}^{(s,k)} = (\mathbf{Y}_1^{(s,k)}, \dots, \mathbf{Y}_h^{(s,k)})$ denote the latent variables imputed for the training sample s at iteration k . Occasionally, we use the notation $\mathbf{Y}_0^{(s,k)} = \mathbf{Y}_0^{(s)} = \mathbf{X}^{(s)}$ and $\mathbf{Y}_{h+1}^{(s,k)} = \mathbf{Y}_{h+1}^{(s)} = \mathbf{Y}^{(s)}$. This algorithm is called “adaptive” as the transition kernel used in step (i) changes with

iterations through the working estimate $\boldsymbol{\theta}^{(k)}$.

S4 Covariance of Latent Variables in the StoNet

Consider the case that we have a regression StoNet trained by the IRO algorithm. In this case, the prediction uncertainty can be quantified by a recursive application of Eve’s law.

Let \mathbf{z} denote a test point at which the prediction uncertainty needs to be quantified. For simplicity of notation, we suppress the bias term by including it as a special column of the corresponding weight matrix. To indicate the iterative nature of the IRO algorithm, we include the superscript ‘ t ’ in the derivation. Let $\mathbf{Z}_i^{(t)}$ denote the imputed latent variable, corresponding to the input \mathbf{z} , for layer i at iteration t . For convenience, we let $\mathbf{Z}_0^{(t)} = \mathbf{z}$ for all t . Let $\boldsymbol{\mu}_i^{(t)}$ and $\boldsymbol{\Sigma}_i^{(t)}$ denote, respectively, the mean and covariance matrix of $\mathbf{Z}_i^{(t)}$. Let $\mathbf{w}_{ij}^{(t)}$ denote the j -th row of the weight matrix $\mathbf{w}_i^{(t)}$, which represents the weights from the neurons of layer $i - 1$ to neuron j of layer i . By Eve’s law, for any layer $i \in \{2, 3, \dots, h + 1\}$, we

Algorithm S2: An adaptive SGHMC algorithm for training StoNet

Input: Dataset (\mathbf{X}, \mathbf{Y}) , total iteration number K , Monte Carlo step number t_{HMC} , the learning rate sequence $\{\epsilon_{k,i} : t = 1, 2, \dots, T; i = 1, 2, \dots, h+1\}$, and the step size sequence $\{\gamma_{k,i} : t = 1, 2, \dots, T; i = 1, 2, \dots, h+1\}$.

Initialization: Randomly initialize the network parameters

$$\hat{\boldsymbol{\theta}}^{(0)} = (\hat{\boldsymbol{\theta}}_1^{(0)}, \dots, \hat{\boldsymbol{\theta}}_{h+1}^{(0)}).$$

for $k = 1, 2, \dots, K$ **do**

STEP 0: Subsampling: Draw a mini-batch of data and denote it by S_k .

STEP 1: Backward Sampling: For each observation $s \in S_k$, sample \mathbf{Y}_i 's in the order from layer h to layer 1. More explicitly, we sample $\mathbf{Y}_i^{(s,k)}$ from the distribution

$$\pi(\mathbf{Y}_i^{(s,k)} | \hat{\boldsymbol{\theta}}_i^{(k-1)}, \hat{\boldsymbol{\theta}}_{i+1}^{(k-1)}, \mathbf{Y}_{i+1}^{(s,k)}, \mathbf{Y}_{i-1}^{(s,k)}) \propto \pi(\mathbf{Y}_{i+1}^{(s,k)} | \hat{\boldsymbol{\theta}}_{i+1}^{(k-1)}, \mathbf{Y}_i^{(s,k)}) \pi(\mathbf{Y}_i^{(s,k)} | \hat{\boldsymbol{\theta}}_i^{(k-1)}, \mathbf{Y}_{i-1}^{(s,k)})$$

 by running SGHMC for t_{HMC} steps:

 Initialize $\mathbf{v}_i^{(s,0)} = \mathbf{0}$, and initialize $\mathbf{Y}_i^{(s,k,0)}$ by forward pass of DNN.

for $l = 1, 2, \dots, t_{HMC}$ **do**

for $i = h, h-1, \dots, 1$ **do**

$$\begin{aligned} \mathbf{v}_i^{(s,k,l)} &= (1 - \epsilon_{k,i} \eta_i) \mathbf{v}_i^{(s,k,l-1)} + \epsilon_{k,i} \nabla_{\mathbf{Y}_i^{(s,k,l-1)}} \log \pi(\mathbf{Y}_i^{(s,k,l-1)} | \hat{\boldsymbol{\theta}}_i^{(k-1)}, \mathbf{Y}_{i-1}^{(s,k,l-1)}) \\ &\quad + \epsilon_{k,i} \nabla_{\mathbf{Y}_i^{(s,k,l-1)}} \log \pi(\mathbf{Y}_{i+1}^{(s,k,l-1)} | \hat{\boldsymbol{\theta}}_{i+1}^{(k-1)}, \mathbf{Y}_i^{(s,k,l-1)}) + \sqrt{2\epsilon_{k,i}\eta} \mathbf{e}^{(s,k,l)}, \\ \mathbf{Y}_i^{(s,k,l)} &= \mathbf{Y}_i^{(s,k,l-1)} + \epsilon_{k,i} \mathbf{v}_i^{(s,k,l-1)}, \end{aligned} \tag{S3.7}$$

 where $\mathbf{e}^{s,k,l} \sim N(0, \mathbf{I}_{d_i})$, $\epsilon_{k,i}$ is the learning rate, and η is the friction coefficient. The algorithm is reduced to SGLD when $\epsilon_{k,i} \eta_i \equiv 1$.

 Set $\mathbf{Y}_i^{(s,k)} = \mathbf{Y}_i^{(s,k,t_{HMC})}$ for $i = 1, 2, \dots, h$.

STEP 2: Parameter Update: Update the estimates of

$\hat{\boldsymbol{\theta}}^{(k-1)} = (\hat{\boldsymbol{\theta}}_1^{(k-1)}, \hat{\boldsymbol{\theta}}_2^{(k-1)}, \dots, \hat{\boldsymbol{\theta}}_{h+1}^{(k-1)})$ by stochastic gradient descent

$$\hat{\boldsymbol{\theta}}_i^{(k)} = \hat{\boldsymbol{\theta}}_i^{(k-1)} + \gamma_{k,i} \left(\frac{n}{|S_k|} \sum_{s \in S_k} \nabla_{\boldsymbol{\theta}_i} \log \pi(\mathbf{Y}_i^{(s,k)} | \hat{\boldsymbol{\theta}}_i^{(k-1)}, \mathbf{Y}_{i-1}^{(s,k)}) - n \nabla_{\boldsymbol{\theta}_i} P_\lambda(\hat{\boldsymbol{\theta}}_i) \right), \tag{S3.8}$$

 for $i = 1, 2, \dots, h+1$, where $\gamma_{k,i}$ is the step size used for updating $\boldsymbol{\theta}_i$.

Output: $\hat{\boldsymbol{\theta}}_n^{(K)}$

then have

$$\begin{aligned}
\boldsymbol{\Sigma}_i^{(t)} &= \mathbb{E}(\text{Var}(\mathbf{Z}_i^{(t)} | \mathbf{Z}_{i-1}^{(t)})) + \text{Var}(\mathbb{E}(\mathbf{Z}_i^{(t)} | \mathbf{Z}_{i-1}^{(t)})) \\
&= \mathbb{E} \text{diag} \left\{ \psi(\mathbf{Z}_{i-1}^{(t)})^T \text{Var}(\hat{\mathbf{w}}_{i_j}^{(t)}) \psi(\mathbf{Z}_{i-1}^{(t)}) : j = 1, \dots, d_i \right\} + \text{Var} \left(\mathbb{E}(\hat{\mathbf{w}}_i) \psi(\mathbf{Z}_{i-1}^{(t)}) \right) \\
&= \text{diag} \left\{ \text{tr}(\text{Var}(\hat{\mathbf{w}}_{i_j}^{(t)}) \text{Var}(\psi(\mathbf{Z}_{i-1}^{(t)}))) + (\mathbb{E}(\psi(\mathbf{Z}_{i-1}^{(t)})))^T \right. \\
&\quad \left. \times \text{Var}(\hat{\mathbf{w}}_{i_j}^{(t)}) (\mathbb{E}(\psi(\mathbf{Z}_{i-1}^{(t)}))) : j = 1, \dots, d_i \right\} + \mathbb{E}(\hat{\mathbf{w}}_i) \text{Var}(\psi(\mathbf{Z}_{i-1}^{(t)})) (\mathbb{E}(\hat{\mathbf{w}}_i))^T,
\end{aligned}$$

where $\text{Var}(\hat{\mathbf{w}}_{i_j}^{(t)})$ is calculated by the Lasso+OLS or Lasso+mLS procedure suggested by Liu and Yu (2013). We refer to Theorem 3 of Liu and Yu (2013) for asymptotic normality of the non-sparse components of $\hat{\mathbf{w}}_{i_j}^{(t)}$. For the OLS case, the non-sparse submatrix of $\text{Var}(\hat{\mathbf{w}}_{i_j}^{(t)})$ is given by

$$\widetilde{\text{Var}(\hat{\mathbf{w}}_{i_j}^{(t)})} = \hat{\varsigma}_{i,j}^2 [(\psi(\tilde{\mathbb{Y}}_{i-1}^{(t)})^T \psi(\tilde{\mathbb{Y}}_{i-1}^{(t)})]^{-1},$$

where $\psi(\tilde{\mathbb{Y}}_{i-1}^{(t)})$ is the design matrix of the linear regression

$$\mathbb{Y}_{i,j}^{(t)} = \psi(\tilde{\mathbb{Y}}_{i-1}^{(t)}) (\tilde{\mathbf{w}}_{i_j}^{(t)})^T + \boldsymbol{\epsilon}_{i,j}$$

selected by Lasso for neuron j of layer i at iteration t , $\boldsymbol{\epsilon}_{i,j} \sim N(0, \varsigma_i^2 I_n)$, and $\hat{\varsigma}_{i,j}^2$ denotes the OLS estimator of ς_i^2 . Here $\mathbb{Y}_{i-1}^{(t)} \in \mathbb{R}^{n \times d_{i-1}}$ denotes imputed latent variables for all neurons of layer $i-1$, $\mathbb{Y}_{i,j}^{(t)} \in \mathbb{R}^n$ denotes imputed latent variables for neuron j of layer i , $\tilde{\mathbb{Y}}_{i-1}^{(t)} \in \mathbb{R}^{n \times \tilde{q}_{i,j}}$ denotes the variables selected by Lasso, $\tilde{\mathbf{w}}_{i_j}^{(t)}$ denotes the corresponding regression coefficients, and $\tilde{q}_{i,j}$ denotes the number of selected variables.

Let $\boldsymbol{\mu}_{i-1}^{(t)} = (\mu_{i-1,1}^{(t)}, \dots, \mu_{i-1,d_{i-1}}^{(t)})^T$ denote the mean of $\mathbf{Z}_{i-1}^{(t)}$, and let $D_{\psi'}(\boldsymbol{\mu}_{i-1}^{(t)}) = \text{diag}\{\psi'(\mu_{i-1,1}^{(t)}), \dots, \psi'(\mu_{i-1,d_{i-1}}^{(t)})\}$, where ψ' denotes the first derivative of the activation function ψ . By the first order Taylor expansion, we have

$$\mathbb{E}(\psi(\mathbf{Z}_{i-1}^{(t)})) \approx \psi(\boldsymbol{\mu}_{i-1}^{(t)}), \quad \text{Var}(\psi(\mathbf{Z}_{i-1}^{(t)})) \approx D_{\psi'}(\boldsymbol{\mu}_{i-1}^{(t)}) \Sigma_{i-1}^{(t)} D_{\psi'}(\boldsymbol{\mu}_{i-1}^{(t)}).$$

Further, if we estimate $\mathbb{E}(\hat{\mathbf{w}}_i)$ by $\hat{\mathbf{w}}_i$ and estimate $\boldsymbol{\mu}_{i-1}^{(t)}$ by $\mathbf{Z}_{i-1}^{(t)}$, then we have the approximation:

$$\begin{aligned} \widehat{\boldsymbol{\Sigma}}_i^{(t)} \approx & \text{diag}\left\{ \text{tr}(\text{Var}(\hat{\mathbf{w}}_{i_j}^{(t)}) D_{\psi'}(\mathbf{Z}_{i-1}^{(t)}) \widehat{\boldsymbol{\Sigma}}_{i-1}^{(t)} D_{\psi'}(\mathbf{Z}_{i-1}^{(t)})) + (\psi(\mathbf{Z}_{i-1}^{(t)}))^T \text{Var}(\hat{\mathbf{w}}_{i_j}^{(t)}) \psi(\mathbf{Z}_{i-1}^{(t)}) : j = 1, \dots, d_i \right\} \\ & + \hat{\mathbf{w}}_i^{(t)} D_{\psi'}(\mathbf{Z}_{i-1}^{(t)}) \widehat{\boldsymbol{\Sigma}}_{i-1}^{(t)} D_{\psi'}(\mathbf{Z}_{i-1}^{(t)}) (\hat{\mathbf{w}}_i^{(t)})^T. \end{aligned} \quad (\text{S4.9})$$

For the first hidden layer, it is reduced to

$$\begin{aligned} \widehat{\boldsymbol{\Sigma}}_1^{(t)} \approx & \text{diag}\left\{ \text{tr}(\text{Var}(\hat{\mathbf{w}}_{1_j}^{(t)}) \text{Var}(\mathbf{z})) + \mathbf{z}^T \text{Var}(\hat{\mathbf{w}}_{1_j}^{(t)}) \mathbf{z} : j = 1, \dots, d_1 \right\} \\ & + \hat{\mathbf{w}}_1^{(t)} \text{Var}(\mathbf{z}) (\hat{\mathbf{w}}_1^{(t)})^T. \end{aligned} \quad (\text{S4.10})$$

Since $\text{Var}(\mathbf{z}) = 0$ holds for any fixed test point \mathbf{z} , $\widehat{\boldsymbol{\Sigma}}_1^{(t)}$ can be further reduced to

$$\widehat{\boldsymbol{\Sigma}}_1^{(t)} \approx \text{diag}\left\{ \mathbf{z}^T \text{Var}(\hat{\mathbf{w}}_{1_j}^{(t)}) \mathbf{z} : j = 1, 2, \dots, d_1 \right\}. \quad (\text{S4.11})$$

S4.1 More Numerical Results

Table S1: Calibration results for CIFAR10 data, where the standard deviations of the respective measures are given in parentheses.

Network	Size	Method	ACC(%)	NLL	ECE(%)	CECE(%)
DenseNet40	176K	No Calibration	92.80 (0.08)	0.3101 (0.0045)	4.45 (0.15)	0.95 (0.03)
		Temp. Scaling	92.80 (0.08)	0.2205 (0.0017)	1.23 (0.09)	0.43 (0.02)
		Matrix Scaling	92.36 (0.15)	0.2277 (0.0034)	1.28 (0.17)	0.41 (0.02)
		Focal	92.04 (0.15)	0.2377 (0.0037)	1.36 (0.09)	0.43 (0.03)
		MMCE	92.24 (0.49)	0.2362 (0.0258)	1.37 (0.35)	0.44 (0.07)
		Post-Linear	92.34 (0.25)	0.2320 (0.0106)	1.04 (0.99)	0.44 (0.17)
		Post-StoNet	92.63 (0.13)	0.2214 (0.0044)	0.54 (0.07)	0.31 (0.05)
ResNet110	1.7M	No Calibration	92.70 (0.90)	0.3359 (0.0472)	4.84 (0.63)	1.02 (0.13)
		Temp. Scaling	92.70 (0.90)	0.2238 (0.0267)	1.29 (0.11)	0.45 (0.03)
		Matrix Scaling	92.15 (0.42)	0.2377 (0.0096)	1.56 (0.14)	0.46 (0.02)
		Focal	91.97 (0.29)	0.2399 (0.0107)	0.87 (0.11)	0.44 (0.03)
		MMCE	91.81 (0.38)	0.2476 (0.0216)	1.57 (0.25)	0.49 (0.07)
		Post-Linear	92.59 (0.44)	0.2230(0.0108)	1.12 (0.29)	0.39 (0.03)
		Post-StoNet	92.66 (0.84)	0.2210 (0.0269)	0.47 (0.23)	0.32 (0.06)
WideResNet-28-10	36M	No Calibration	95.84 (0.18)	0.1704 (0.0037)	2.55 (0.09)	0.56 (0.01)
		Temp. Scaling	95.84 (0.18)	0.1468 (0.0029)	1.16 (0.05)	0.34 (0.02)
		Matrix Scaling	93.67 (0.81)	0.1961 (0.0218)	1.47 (0.18)	0.41 (0.02)
		Focal	95.43 (0.07)	0.1943 (0.0149)	6.29 (1.33)	1.37 (0.28)
		MMCE	93.68 (0.81)	0.1993 (0.0236)	1.43 (0.09)	0.46 (0.03)
		Post-Linear	95.77 (0.13)	0.1444 (0.0060)	0.94 (0.40)	0.30 (0.10)
		Post-StoNet	95.64 (0.12)	0.1449 (0.0018)	0.87 (0.11)	0.25 (0.02)

S5 Hyper-parameter Settings for the Numerical Experiments

For Algorithm S2, since the learning rates $\epsilon_{k,i}$'s and the latent variable variances σ_i^2 's can be largely canceled at each step of latent variable imputation, their absolute values do not mean much to the convergence of the simulation. For this reason, we often set their values to be very small in our numerical experiments, which merely controls the size of random noise

added to the corresponding latent variables.

S5.1 Settings for the Illustrative Example

One-hidden-layer StoNet: we tried three parameter settings:

- (i) $\sigma_2^2 = 5e-5, \sigma_1^2 = 5e-6, \epsilon_{k,1} = 5e-9, \eta_i = \frac{1}{\epsilon_{k,i}}, t_{HMC} = 1, \frac{\gamma_{k,1}}{|S_k|} = 5e-4,$
 $\frac{\gamma_{k,2}}{|S_k|} = 5e-8, |S_k| = 50;$
- (ii) $\sigma_2^2 = 1e-4, \sigma_1^2 = 1e-5, \epsilon_{k,1} = 1e-8, \eta_i = \frac{1}{\epsilon_{k,i}}, t_{HMC} = 1, \frac{\gamma_{k,1}}{|S_k|} = 5e-4,$
 $\frac{\gamma_{k,2}}{|S_k|} = 5e-8, |S_k| = 50;$
- (iii) $\sigma_2^2 = 2e-4, \sigma_1^2 = 2e-5, \epsilon_{k,1} = 2e-8, \eta_i = \frac{1}{\epsilon_{k,i}}, t_{HMC} = 1, \frac{\gamma_{k,1}}{|S_k|} = 5e-4,$
 $\frac{\gamma_{k,2}}{|S_k|} = 5e-8, |S_k| = 50.$

Two-hidden-layer StoNet: we tried three parameter settings:

- (i) $\sigma_3^2 = 5e-10, \sigma_2^2 = 5e-11, \sigma_1^2 = 5e-12, \epsilon_{k,2} = 5e-14, \epsilon_{k,1} = 1e-14,$
 $\eta_i = \frac{1}{\epsilon_{k,i}}, t_{HMC} = 1, \frac{\gamma_{k,3}}{|S_k|} = 5e-6, \frac{\gamma_{k,2}}{|S_k|} = 5e-10, \frac{\gamma_{k,1}}{|S_k|} = 5e-14,$
 $|S_k| = 50;$
- (ii) $\sigma_3^2 = 1e-9, \sigma_2^2 = 1e-10, \sigma_1^2 = 1e-11, \epsilon_{k,2} = 1e-13, \epsilon_{k,1} = 1e-14,$
 $\eta_i = \frac{1}{\epsilon_{k,i}}, t_{HMC} = 1, \frac{\gamma_{k,3}}{|S_k|} = 5e-6, \frac{\gamma_{k,2}}{|S_k|} = 5e-10, \frac{\gamma_{k,1}}{|S_k|} = 5e-14,$
 $|S_k| = 50;$
- (iii) $\sigma_3^2 = 2e-9, \sigma_2^2 = 2e-10, \sigma_1^2 = 2e-11, \epsilon_{k,2} = 1e-13, \epsilon_{k,1} = 1e-14,$
 $\eta_i = \frac{1}{\epsilon_{k,i}}, t_{HMC} = 1, \frac{\gamma_{k,3}}{|S_k|} = 5e-6, \frac{\gamma_{k,2}}{|S_k|} = 5e-10, \frac{\gamma_{k,1}}{|S_k|} = 5e-14, |S_k| = 50.$

For both StoNets, the major difference among the settings is at σ_i 's. For convenience, we call the settings (i), (ii) and (iii), respectively.

S5.2 Settings for the Experiments in Section 5

CIFAR100 and CIFAR10: Following the setting of post-calibration methods in Guo et al. (2017), we split the training data into a training set of 45,000 images and a holdout validation set of 5,000 images. The training settings for the three models are:

- *ResNet110:* The model was trained on the training set using SGD with momentum for 200 epochs with the batch size 128, momentum 0.9, and weight decay 0.0001. The learning rate was set to 0.1 for the first 80 epochs and divided by 10 at the 80-th and 150-th epochs.
- *Densenet40:* The model was trained on the training set using SGD with momentum for 300 epochs with the batch size 128, momentum 0.9, and weight decay 0.0001. The learning rate was set to 0.1 for the first 150 epochs and divided by 10 at the 150-th and 225-th epochs.
- *WideResNet-28-10:* The model was trained on the training set using SGD with momentum for 200 epochs with momentum 0.9, and weight decay 0.0005. The learning rate was set to 0.1 for the first 60 epochs and divided by 10 at 60-th, 120-th, and 160-th epochs.

After training, we extracted the outputs of the last fully connected layer of each model on the validation set, and used them as input to a StoNet model with one hidden layer, 100 hidden units, and the activation function *tanh*. The StoNet model was trained using Algorithm (S2) with the hyper-parameters as given in Table S2. The regularization parameter λ was set to $1e-4$ for CIFAR10 and $5e-5$ for CIFAR100. As a baseline, we consider the Post-Linear model. We used the outputs of the last fully connected layer of each model as input features, and trained sparse multi-class logistic regression models with LASSO penalty on the validation set. The regularization parameter is selected via a 5-fold cross-validation using the default setting in the *scikit-learn* package.

Table S2: Post-StoNet Hyper-Parameter Setting for CIFAR10 and CIFAR100 data

Hyper-Parameter	Value
$[\sigma_1^2, \sigma_2^2]$	$[1e-2, 1e-3]$
$\epsilon_{k,1}$	$1e-7$
η_1	$\frac{1}{\epsilon_{k,1}}$
t_{HMC}	1
$[\gamma_{k,1}, \gamma_{k,2}]$	$[\frac{5e-4}{5000}, \frac{5e-6}{5000}]$
$ S_k $	50
$P_\lambda(\boldsymbol{\theta})$	$\lambda \ \boldsymbol{\theta}\ _1$

Regression Examples: The data sets were from UCI machine learning repository. For all experiments, we split the data into 40% as the training set, 40% as the calibration/validation set (used to fit a StoNet model for

our approach and to compute the absolute value of the residue as the non-conformity score for Split Conformal), and 20% as the test set. The random split was repeated 10 times. We reported the mean values and standard deviations of the prediction interval length and coverage rate.

We modeled each dataset using a DNN with 2 hidden layers, with 1000 and 100 hidden units respectively, and the activation function *tanh*. The DNN was trained using Adam (Kingma and Ba, 2015) with a batch size of 50 and a constant learning rate of 0.001. The algorithm was run for 1000 epochs for the Protein data set, 200 epochs for the Year data set, and 5000 epochs for other data sets. After the DNN was trained, we refit a StoNet on the calibration set using the output of the last hidden layer of the DNN as input. The StoNet had one hidden layer, 20 hidden units, and the activation function *tanh*. Algorithm S2 was used to train the StoNet with the hyperparameters as given in Table S3. The penalty parameter λ was selected from $\{1e-1, 8e-2, 5e-2, 4e-2, 3e-2, 2e-2, 1e-2, 5e-3, 2e-3, 1e-3, 5e-4\}$ by 5-fold cross-validation, where we picked a value of λ such that the average coverage rate on the calibration sets were closest to the target level 90%. Specifically, we picked $\lambda = 8e-2$ for the Liver dataset, $\lambda = 8e-2$ for QSAR dataset, $\lambda = 3e-3$ for Community dataset, $\lambda = 8e-2$ for STAR

<https://archive.ics.uci.edu/dataset/60/liver+disorders>
<https://archive.ics.uci.edu/dataset/504/qsar+fish+toxicity>
<https://archive.ics.uci.edu/dataset/183/communities+and+crime>

dataset, $\lambda = 5e - 2$ for Abalone dataset, $\lambda = 5e - 2$ for Parkinson dataset, $\lambda = 5e - 3$ for Power Plant dataset, $\lambda = 2e - 3$ for Bike dataset, $\lambda = 5e - 3$ for Protein dataset, $\lambda = 1e - 3$ for Year dataset

Table S3: StoNet Hyper-Parameter Setting for UCI data sets, where N is size of the calibration data set.

Hyper-Parameter	Value
$[\sigma_1^2, \sigma_2^2]$	$[1e-4, 1e-5]$
$\epsilon_{k,1}$	$1e-7$
η_1	$\frac{1}{\epsilon_{k,1}}$
t_{HMC}	1
$[\gamma_{k,1}, \gamma_{k,2}]$	$[\frac{1e-3}{N}, \frac{1e-5}{N}]$
$ S_k $	50
$P_\lambda(\boldsymbol{\theta})$	$\lambda \ \boldsymbol{\theta}\ _1$

S6 Consistency is Essential for the Validity of the Post-StoNet Approach

The parameter estimation consistency is essential for the validity of the post-StoNet approach. To demonstrate this issue, we applied the post-StoNet modeling approach to the Community dataset without regularization (i.e. setting $\lambda = 0$), which violates the sparsity condition of Theorem

<https://github.com/yromano/cqr/tree/master/datasets>
<https://archive.ics.uci.edu/dataset/1/abalone>
<https://archive.ics.uci.edu/dataset/189/parkinsons+telemonitoring>
<https://archive.ics.uci.edu/dataset/294/combined+cycle+power+plant>
<https://archive.ics.uci.edu/ml/datasets/bike+sharing+dataset>
<https://archive.ics.uci.edu/dataset/265/physicochemical+properties+of+protein+tertiary+structure>
<https://archive.ics.uci.edu/dataset/203/yearpredictionmsd>

1. The resulting prediction intervals have only a coverage rate of 29.25% (with a standard deviation of 3.77%).

Bibliography

Chen, T., E. Fox, and C. Guestrin (2014). Stochastic gradient hamiltonian monte carlo. In *International conference on machine learning*, pp. 1683–1691.

Deng, W., X. Zhang, F. Liang, and G. Lin (2019). An adaptive empirical bayesian method for sparse deep learning. *Advances in neural information processing systems 2019*, 5563–5573.

Guo, C., G. Pleiss, Y. Sun, and K. Q. Weinberger (2017). On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR.

Huang, J., S. Ma, and C.-H. Zhang (2008). The iterated lasso for high-dimensional logistic regression. *The University of Iowa, Department of Statistics and Actuarial Sciences*.

Kingma, D. and J. Ba (2015). Adam: a method for stochastic optimization. In *International Conference on Learning Representations*.

- Liang, F., B. Jia, J. Xue, Q. Li, and Y. Luo (2018). An imputation-regularized optimization algorithm for high dimensional missing data problems and beyond. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80, 899–926.
- Liang, S., Y. Sun, and F. Liang (2022). Nonlinear sufficient dimension reduction with a stochastic neural network. *NeurIPS* 35.
- Liu, H. and B. Yu (2013). Asymptotic properties of Lasso+mLS and Lasso+Ridge in sparse high-dimensional linear regression. *Electronic Journal of Statistics* 7, 3124 – 3169.
- Meinshausen, N. and B. Yu (2009). Lasso-type recovery of sparse representations for high-dimensional data. *Annals of Statistics* 37, 246–270.