

---

# Sparse Deep Learning: A New Framework Immune to Local Traps and Miscalibration

---

**Yan Sun**

Purdue University  
West Lafayette, IN 47906  
sun748@purdue.edu

**Wenjun Xiong**

Guangxi Normal University & Purdue University  
West Lafayette, IN 47906  
xiong90@purdue.edu

**Faming Liang**

Purdue University  
West Lafayette, IN 47906  
fmliang@purdue.edu

## 1 Proof of Theorem 2.1

*Proof.* We first define the equivalent class of neural network parameters. Given a parameter vector  $\beta$  and the corresponding structure parameter vector  $\gamma$ , its equivalent class is given by

$$Q_E(\beta, \gamma) = \{(\tilde{\beta}, \tilde{\gamma}) : \nu_g(\tilde{\beta}, \tilde{\gamma}) = (\beta, \gamma), \mu(\tilde{\beta}, \tilde{\gamma}, \mathbf{x}) = \mu(\beta, \gamma, \mathbf{x}), \forall \mathbf{x}\},$$

where  $\nu_g(\cdot)$  denotes a generic mapping that contains only the transformations of node permutation and weight sign flipping. Specifically, we define

$$Q_E^* = Q_E(\beta^*, \gamma^*),$$

which represents the equivalent class of *true DNN model*.

Let  $B_{\delta_n}(\beta^*) = \{\beta : |\beta_i - \beta_i^*| < \delta_n, \forall i \in \gamma^*, |\beta_i - \beta_i^*| < 2\sigma_{0,n} \log(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}}), \forall i \notin \gamma^*\}$ . By assumption C.1,  $\beta^*$  is generic (i.e.  $Q_E(\beta^*)$  contains only reparameterizations of weight sign-flipping or node permutations as defined in Feng and Simon (2017) and Fefferman (1994)) and  $\min_{i \in \gamma^*} |\beta_i^*| - \delta_n > (C-1)\delta_n > \delta_n$ , then for any  $\beta^{*(1)}, \beta^{*(2)} \in Q_E^*$ ,  $B_{\delta_n}(\beta^{*(1)}) \cap B_{\delta_n}(\beta^{*(2)}) = \emptyset$ , and thus  $\{\beta : \tilde{\nu}(\beta) \in B_{\delta_n}(\beta^*)\} = \cup_{\beta \in Q_E^*} B_{\delta_n}(\beta)$ . In what follows, we will first show  $\pi(\cup_{\beta \in Q_E^*} B_{\delta_n}(\beta) | D_n) \rightarrow 1$  as  $n \rightarrow \infty$ , which means the most posterior mass falls in the neighbourhood of true parameter.

Remark on the notation:  $\tilde{\nu}(\cdot)$  is similar to  $\nu(\cdot)$  defined in Section 2.1 of the main text. They both map the set  $Q_E(\beta, \gamma)$  to a unique network. The difference between them is that  $\|\nu(\beta) - \beta^*\|_\infty$  may be arbitrary, but  $\|\tilde{\nu}(\beta) - \beta^*\|_\infty$  is minimized. In other words,  $\nu(\beta, \gamma)$  is to find an arbitrary network in  $Q_E(\beta, \gamma)$  as the representative of the equivalent class, while  $\tilde{\nu}(\beta, \gamma)$  is to find a representative in  $Q_E(\beta, \gamma)$  such that the distance between  $\beta^*$  and the representative is minimized. In what follows, we will use  $\tilde{\nu}(\beta)$  and  $\tilde{\nu}(\gamma)$  to denote the connection weight and network structure of  $\tilde{\nu}(\beta, \gamma)$ , respectively. With a slight abuse of notation, we will use  $\tilde{\nu}(\beta)_i$  to denote the  $i$ th component of  $\tilde{\nu}(\beta)$ , and use  $\tilde{\nu}(\gamma)_i$  to denote the  $i$ th component of  $\tilde{\nu}(\gamma)$ .

Recall that the marginal posterior inclusion probability is given by

$$q_i = \int \sum_{\gamma} e_{i|\tilde{\nu}(\beta, \gamma)} \pi(\gamma | \beta, D_n) \pi(\beta | D_n) d\beta = \int \pi(\tilde{\nu}(\gamma)_i = 1 | \beta) \pi(\beta | D_n) d\beta.$$

For the mixture Gaussian prior,

$$\pi(\gamma_i = 1 | \beta) = \frac{1}{1 + \frac{\sigma_{1,n}(1-\lambda_n)}{\sigma_{0,n}\lambda_n} e^{-\left(\frac{1}{2\sigma_{0,n}^2} - \frac{1}{2\sigma_{1,n}^2}\right)\beta_i^2}},$$

which increases with respect to  $|\beta_i|$ . In particular, if  $|\beta_i| > 2\sigma_{0,n} \log(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}})$ , then  $\pi(\gamma_i = 1 | \beta) > \frac{1}{2}$ .

For the mixture Gaussian prior,

$$\begin{aligned} & \pi(\beta \notin \cup_{\beta \in Q_E^*} B_{\delta_n}(\beta) \mid D_n) \\ & \leq \pi(\exists i \notin \gamma^*, |\tilde{\nu}(\beta)_i| > 2\sigma_{0,n} \log(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}}) \mid D_n) + \pi(\exists i \in \gamma^*, |\tilde{\nu}(\beta)_i - \beta_i^*| > \delta_n \mid D_n). \end{aligned}$$

For the first term, note that for a given  $i \notin \gamma^*$ ,

$$\begin{aligned} \pi(|\tilde{\nu}(\beta)_i| > 2\sigma_{0,n} \log(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}}) \mid D_n) & \leq \pi(\pi(\tilde{\nu}(\gamma)_i = 1 | \beta) > \frac{1}{2} \mid D_n) \\ & \leq 2 \int \pi(\tilde{\nu}(\gamma)_i = 1 | \beta) \pi(\beta | D_n) d\beta \\ & \leq 2\rho(\epsilon_n) + 2\pi(d(p_\beta, p_{\mu^*}) \geq \epsilon_n \mid D_n) \rightarrow 0. \end{aligned}$$

Then we have

$$\begin{aligned} \pi(\exists i \notin \gamma^*, |\tilde{\nu}(\beta)_i| > 2\sigma_{0,n} \log(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}}) \mid D_n) & = \pi(\max_{i \notin \gamma^*} |\tilde{\nu}(\beta)_i| > 2\sigma_{0,n} \log(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}}) \mid D_n) \\ & \leq \pi(\max_{i \notin \gamma^*} \pi(\tilde{\nu}(\gamma)_i = 1 | \beta) > \frac{1}{2} \mid D_n) \\ & \leq \sum_{i \notin \gamma^*} \pi(\pi(\tilde{\nu}(\gamma)_i = 1 | \beta) > \frac{1}{2} \mid D_n) \\ & \leq 2K_n \rho(\epsilon_n) + 2K_n \pi(d(p_\beta, p_{\mu^*}) \geq \epsilon_n \mid D_n) \rightarrow 0. \end{aligned}$$

For the second term, by condition C.1 and Lemma 2.1,

$$\begin{aligned} & \pi(\exists i \in \gamma^*, |\tilde{\nu}(\beta)_i - \beta_i^*| > \delta_n \mid D_n) = \pi(\max_{i \in \gamma^*} |\tilde{\nu}(\beta)_i - \beta_i^*| > \delta_n \mid D_n) \\ & = \pi(\max_{i \in \gamma^*} |\tilde{\nu}(\beta)_i - \beta_i^*| > \delta_n, d(p_\beta, p_{\mu^*}) \leq \epsilon_n \mid D_n) \\ & \quad + \pi(\max_{i \in \gamma^*} |\tilde{\nu}(\beta)_i - \beta_i^*| > \delta_n, d(p_\beta, p_{\mu^*}) \geq \epsilon_n \mid D_n) \\ & \leq \pi(A(\epsilon_n, \delta_n) \mid D_n) + \pi(d(p_\beta, p_{\mu^*}) \geq \epsilon_n \mid D_n) \rightarrow 0. \end{aligned}$$

Summarizing the above two terms, we have  $\pi(\cup_{\beta \in Q_E^*} B_{\delta_n}(\beta) \mid D_n) \rightarrow 1$ .

Let  $Q_n = |Q_E^*|$  be the number of equivalent *true DNN model*. By the generic assumption of  $\beta^*$ , for any  $\beta^{*(1)}, \beta^{*(2)} \in Q_E^*$ ,  $B_{\delta_n}(\beta^{*(1)}) \cap B_{\delta_n}(\beta^{*(2)}) = \emptyset$ . Then in  $B_{\delta_n}(\beta^*)$ , the posterior density of  $\tilde{\nu}(\beta)$  is  $Q_n$  times the posterior density of  $\beta$ . Then for a function  $f(\cdot)$  of  $\tilde{\nu}(\beta)$ , by changing variable,

$$\int_{\tilde{\nu}(\beta) \in B_{\delta_n}(\beta^*)} f(\tilde{\nu}(\beta)) \pi(\tilde{\nu}(\beta) | D_n) d\tilde{\nu}(\beta) = Q_n \int_{B_{\delta_n}(\beta^*)} f(\beta) \pi(\beta | D_n) d\beta.$$

Thus, we only need to consider the integration on  $B_{\delta_n}(\beta^*)$ . Define

$$\hat{\beta}_i = \begin{cases} \beta_i^* - \sum_{j \in \gamma^*} h^{i,j}(\beta^*) h_j(\beta^*), & \forall i \in \gamma^*, \\ 0, & \forall i \notin \gamma^*. \end{cases}$$

We will first prove that for any real vector  $\mathbf{t}$ ,

$$\begin{aligned} E(e^{\sqrt{n}\mathbf{t}^T(\tilde{\nu}(\beta) - \hat{\beta})} \mid D_n, B_{\delta_n}(\beta^*)) & := \frac{\int_{B_{\delta_n}(\beta^*)} e^{\sqrt{n}\mathbf{t}^T(\tilde{\nu}(\beta) - \hat{\beta})} \pi(\tilde{\nu}(\beta) | D_n) d\tilde{\nu}(\beta)}{\int_{B_{\delta_n}(\beta^*)} \pi(\tilde{\nu}(\beta) | D_n) d\tilde{\nu}(\beta)} \\ & = \frac{\int_{B_{\delta_n}(\beta^*)} e^{\sqrt{n}\mathbf{t}^T(\beta - \hat{\beta})} e^{n\mathbf{l}_n(\beta)} \pi(\beta) d\beta}{\int_{B_{\delta_n}(\beta^*)} e^{n\mathbf{l}_n(\beta)} \pi(\beta) d\beta} \quad (1) \\ & = e^{\frac{1}{2}\mathbf{t}^T \mathbf{V} \mathbf{t} + o_{P^*}(1)}. \end{aligned}$$

For any  $\beta \in B_{\delta_n}(\beta^*)$ , we have

$$\begin{aligned} |\sqrt{n}(\mathbf{t}^T(\beta - \beta_{\gamma^*}))| &\leq \sqrt{n}K_n \|\mathbf{t}\|_{\infty} 2\sigma_{0,n} \log\left(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}}\right) = o(1), \\ |n(l_n(\beta) - l_n(\beta_{\gamma^*}))| &= |n \sum_{i \notin \gamma^*} \beta_i(h_i(\tilde{\beta}))| \leq nK_n M 2\sigma_{0,n} \log\left(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}}\right) = o(1). \end{aligned}$$

Then, we have

$$\begin{aligned} \sqrt{n}\mathbf{t}^T(\beta - \hat{\beta}) &= \sqrt{n}\mathbf{t}^T(\beta - \beta_{\gamma^*} + \beta_{\gamma^*} - \beta^*) + \sqrt{n} \sum_{i,j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_j h_i(\beta^*) \\ &= o(1) + \sqrt{n} \sum_{i \in \gamma^*} (\beta_i - \beta_i^*) \mathbf{t}_i + \sqrt{n} \sum_{i,j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_j h_i(\beta^*), \\ nl_n(\beta) - nl_n(\beta^*) &= n(l_n(\beta) - l_n(\beta_{\gamma^*}) + l_n(\beta_{\gamma^*}) - nl_n(\beta^*)) \\ &= o(1) + n \sum_{i \in \gamma^*} (\beta_i - \beta_i^*) h_i(\beta^*) + \frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) (\beta_i - \beta_i^*) (\beta_j - \beta_j^*) \\ &\quad + \frac{n}{6} \sum_{i,j,k \in \gamma^*} h_{i,j,k}(\tilde{\beta}) (\beta_i - \beta_i^*) (\beta_j - \beta_j^*) (\beta_k - \beta_k^*), \end{aligned} \tag{2}$$

where  $\tilde{\beta}$  is a point between  $\beta_{\gamma^*}$  and  $\beta^*$ . Note that for  $\beta \in B_{\delta_n}(\beta^*)$ ,  $|\beta_i - \beta_i^*| \leq \delta_n \lesssim \frac{1}{\sqrt[3]{nr_n}}$ , we have  $\frac{n}{6} \sum_{i,j,k \in \gamma^*} h_{i,j,k}(\tilde{\beta}) (\beta_i - \beta_i^*) (\beta_j - \beta_j^*) (\beta_k - \beta_k^*) = o(1)$ .

Let  $\beta^{(t)}$  be network parameters satisfying  $\beta_i^{(t)} = \beta_i + \frac{1}{\sqrt{n}} \sum_{j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_j$ ,  $\forall i \in \gamma^*$  and  $\beta_i^{(t)} = \beta_i$ ,  $\forall i \notin \gamma^*$ . Note that  $\frac{1}{\sqrt{n}} \sum_{j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_j \leq \frac{r_n \|\mathbf{t}\|_{\infty} M}{\sqrt{n}} \lesssim \delta_n$ , for large enough  $n$ ,  $|\beta_i^{(t)}| < 2\delta_n$   $\forall i \in \gamma^*$ . Thus, we have

$$\begin{aligned} nl_n(\beta^{(t)}) - nl_n(\beta^*) &= n(l_n(\beta^{(t)}) - l_n(\beta_{\gamma^*}^{(t)}) + l_n(\beta_{\gamma^*}^{(t)}) - nl_n(\beta^*)) \\ &= o(1) + n \sum_{i \in \gamma^*} (\beta_i^{(t)} - \beta_i^*) h_i(\beta^*) + \frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) (\beta_i^{(t)} - \beta_i^*) (\beta_j^{(t)} - \beta_j^*) \\ &= o(1) + n \sum_{i \in \gamma^*} (\beta_i - \beta_i^*) h_i(\beta^*) + \frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) (\beta_i - \beta_i^*) (\beta_j - \beta_j^*) \\ &\quad + \sqrt{n} \sum_{i,j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_j h_i(\beta^*) + \sqrt{n} \sum_{i \in \gamma^*} (\beta_i - \beta_i^*) \mathbf{t}_i + \frac{1}{2} \sum_{i,j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_i \mathbf{t}_j \\ &= o(1) + \sqrt{n}\mathbf{t}^T(\beta - \hat{\beta}) + nl_n(\beta) - nl_n(\beta^*) + \frac{1}{2} \sum_{i,j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_i \mathbf{t}_j, \end{aligned} \tag{3}$$

where the last equality is derived by replacing appropriate terms by  $\sqrt{n}\mathbf{t}^T(\beta - \hat{\beta})$  and  $nl_n(\beta) - nl_n(\beta^*)$  based on (2) and (3), respectively; and the third equality is derived based on the following calculation:

$$\begin{aligned} &\frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) (\beta_i^{(t)} - \beta_i^*) (\beta_j^{(t)} - \beta_j^*) \\ &= \frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) (\beta_i - \beta_i^* + \frac{1}{\sqrt{n}} \sum_{k \in \gamma^*} h^{i,k}(\beta^*) \mathbf{t}_k) (\beta_j - \beta_j^* + \frac{1}{\sqrt{n}} \sum_{k \in \gamma^*} h^{j,k}(\beta^*) \mathbf{t}_k) \\ &= \frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) (\beta_i - \beta_i^*) (\beta_j - \beta_j^*) + 2 \times \frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) \frac{1}{\sqrt{n}} \sum_{k \in \gamma^*} h^{i,k}(\beta^*) \mathbf{t}_k (\beta_j - \beta_j^*) \\ &\quad + \frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) (\frac{1}{\sqrt{n}} \sum_{k \in \gamma^*} h^{i,k}(\beta^*) \mathbf{t}_k) (\frac{1}{\sqrt{n}} \sum_{k \in \gamma^*} h^{j,k}(\beta^*) \mathbf{t}_k) \\ &= \frac{n}{2} \sum_{i,j \in \gamma^*} h_{i,j}(\beta^*) (\beta_i - \beta_i^*) (\beta_j - \beta_j^*) + \sqrt{n} \sum_{i \in \gamma^*} (\beta_i - \beta_i^*) \mathbf{t}_i + \frac{1}{2} \sum_{i,j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_i \mathbf{t}_j, \end{aligned} \tag{4}$$

(5)

where the second and third terms in the last equality are derived based on the relation  $\sum_{i \in \gamma^*} h_{i,j}(\beta^*) h_{i,k}(\beta^*) = \delta_{j,k}$ , where  $\delta_{j,k} = 1$  if  $j = k$ ,  $\delta_{j,k} = 0$  if  $j \neq k$ .

By rearranging the terms in (4), we have

$$\begin{aligned} & \int_{B_{\delta_n}(\beta^*)} \exp\{\sqrt{n}\mathbf{t}^T(\beta - \hat{\beta}) + nl_n(\beta)\} \pi(\beta) d\beta \\ &= \exp\left\{-\frac{1}{2} \sum_{i,j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_i \mathbf{t}_j + o(1)\right\} \int_{B_{\delta_n}(\beta^*)} e^{nl_n(\beta^{(t)})} \pi(\beta) d\beta. \end{aligned}$$

For  $\beta \in B_{\delta_n}(\beta^*)$ ,  $i \in \gamma^*$ , by Assumption C.1, there exists a constant  $C > 2$  such that

$$\begin{aligned} |\beta_i^{(t)}| &\geq |\beta_i| - \frac{r_n \|\mathbf{t}\|_\infty M}{\sqrt{n}} \geq |\beta_i^*| - 2\delta_n \geq (C-2)\delta_n \gtrsim \frac{r_n}{\sqrt{n}} \\ &\gtrsim \sqrt{\left(\frac{1}{2\sigma_{0,n}^2} - \frac{1}{2\sigma_{1,n}^2}\right)^{-1} \log\left(\frac{r_n(1-\lambda_n)\sigma_{1,n}}{\sigma_{0,n}\lambda_n}\right)}. \end{aligned}$$

Then we have

$$\frac{\sigma_{1,n}(1-\lambda_n)}{\sigma_{0,n}\lambda_n} e^{-\left(\frac{1}{2\sigma_{0,n}^2} - \frac{1}{2\sigma_{1,n}^2}\right)(\beta_i^{(t)})^2} \lesssim \frac{1}{r_n}.$$

It is easy to see that the above formula also holds if we replace  $\beta_i^{(t)}$  by  $\beta_i$ . Note that the mixture Gaussian prior of  $\beta_i$  can be written as

$$\pi(\beta_i) = \frac{\lambda_n}{\sqrt{2\pi}\sigma_{1,n}} e^{-\frac{\beta_i^2}{2\sigma_{1,n}^2}} \left(1 + \frac{\sigma_{1,n}(1-\lambda_n)}{\sigma_{0,n}\lambda_n} e^{-\left(\frac{1}{2\sigma_{0,n}^2} - \frac{1}{2\sigma_{1,n}^2}\right)\beta_i^2}\right).$$

Since  $|\beta_i - \beta_i^{(t)}| \lesssim \delta_n \lesssim \frac{1}{\sqrt[3]{nr_n}}$ ,  $|\beta_i + \beta_i^{(t)}| < 2E_n + 3\delta_n \lesssim E_n$ , and  $\frac{1}{\sigma_{1,n}^2} \lesssim \frac{H_n \log(n) + \log(\bar{L})}{E_n^2}$ , we have

$$\frac{r_n}{\sigma_{1,n}^2} (\beta_i - \beta_i^{(t)})(\beta_i + \beta_i^{(t)}) = \frac{H_n \log(n) + \log(\bar{L})}{n^{C_1+1/3}} = o(1),$$

by the condition  $C_1 > 2/3$  and  $H_n \log(n) + \log(\bar{L}) \prec n^{1-\epsilon}$ . Thus,  $\frac{\pi(\beta)}{\pi(\beta^{(t)})} = \prod_{i \in \gamma^*} \frac{\pi(\beta_i)}{\pi(\beta_i^{(t)})} = 1 + o(1)$ , and

$$\begin{aligned} \int_{B_{\delta_n}(\beta^*)} e^{nl_n(\beta^{(t)})} \pi(\beta) d\beta &= (1 + o(1)) \int_{\beta^{(t)} \in B_{\delta_n}(\beta^*)} e^{nl_n(\beta^{(t)})} \pi(\beta^{(t)}) d\beta^{(t)} \\ &= (1 + o(1)) C_N \pi(\beta^{(t)} \in B_{\delta_n}(\beta^*) \mid D_n), \end{aligned} \quad (6)$$

where  $C_N$  is the normalizing constant of the posterior. Note that  $\|\beta^{(t)} - \beta\|_\infty \lesssim \delta_n$ , we have  $\pi(\beta^{(t)} \in B_{\delta_n}(\beta^*) \mid D_n) \rightarrow \pi(\beta \in B_{\delta_n}(\beta^*) \mid D_n)$ . Moreover, since  $-\frac{1}{2} \sum_{i,j \in \gamma^*} h^{i,j}(\beta^*) \mathbf{t}_i \mathbf{t}_j \rightarrow \frac{1}{2} \mathbf{t}^T \mathbf{V} \mathbf{t}$ , we have

$$E(e^{\sqrt{n}\mathbf{t}^T(\tilde{\nu}(\beta) - \hat{\beta})} \mid D_n, B_{\delta_n}(\beta^*)) = \frac{\int_{B_{\delta_n}(\beta^*)} e^{\sqrt{n}\mathbf{t}^T(\beta - \hat{\beta})} e^{nh_n(\beta)} \pi(\beta) d\beta}{\int_{B_{\delta_n}(\beta^*)} e^{nh_n(\beta)} \pi(\beta) d\beta} = e^{\frac{\mathbf{t}^T \mathbf{V} \mathbf{t}}{2} + o_P(1)}.$$

Combining the above result with the fact that  $\pi(\tilde{\nu}(\beta) \in B_{\delta_n}(\beta^*) \mid D_n) \rightarrow 1$ , by section 1 of Castillo and Rousseau (2015), we have

$$\pi[\sqrt{n}(\tilde{\nu}(\beta) - \hat{\beta}) \mid D_n] \rightsquigarrow N(0, \mathbf{V}).$$

We will then show that  $\hat{\beta}$  will converge to  $\beta^*$ , then essentially we can replace  $\hat{\beta}$  by  $\beta^*$  in the above result. Let  $\Theta_{\gamma^*} = \{\beta : \beta_i = 0, \forall i \notin \gamma^*\}$  be the parameter space given the model  $\gamma^*$ , and let  $\hat{\beta}_{\gamma^*}$  be the maximum likelihood estimator given the model  $\gamma^*$ , i.e.

$$\hat{\beta}_{\gamma^*} = \arg \max_{\beta \in \Theta_{\gamma^*}} l_n(\beta).$$

Given condition C.3 and by Theorem 2.1 of Portnoy (1988), we have  $\|\hat{\beta}_{\gamma^*} - \beta^*\| = O(\sqrt{\frac{r_n}{n}}) = o(1)$ . Note that  $h_i(\hat{\beta}_{\gamma^*}) = 0$  as  $\hat{\beta}_{\gamma^*}$  is maximum likelihood estimator. Then for any  $i \in \gamma^*$ ,  $|h_i(\beta^*)| = |h_i(\hat{\beta}_{\gamma^*}) - h_i(\beta^*)| = |\sum_{j \in \gamma^*} h_{ij}(\tilde{\beta})((\hat{\beta}_{\gamma^*})_j - \beta_j^*)| \leq M \|\hat{\beta}_{\gamma^*} - \beta^*\|_1 = O(\sqrt{\frac{r_n}{n}})$ .

Then for any  $i, j \in \gamma^*$ , we have  $\sum_{j \in \gamma^*} h^{i,j}(\beta^*) h_j(\beta^*) = O(\sqrt{\frac{r_n^3}{n}}) = o(1)$ . By the definition of  $\hat{\beta}$ , we have  $\hat{\beta} - \beta^* = o(1)$ . Therefore, we have

$$\pi[\sqrt{n}(\tilde{\nu}(\beta) - \beta^*) \mid D_n] \rightsquigarrow N(0, \mathbf{V}).$$

## 2 Proof of Theorem 2.2

*Proof.* The proof of Theorem 2.2 can be done using the same strategy as that used in proving Theorem 2.1. Here we provide a simpler proof using the result of Theorem 2.1. The notations we used in this proof are the same as in the proof of Theorem 2.1. In the proof of Theorem 2.1, we have shown that  $\pi(\tilde{\nu}(\beta) \in B_{\delta_n}(\beta^*) \mid D_n) \rightarrow 1$ . Note that  $\mu(\beta, \mathbf{x}_0) = \mu(\tilde{\nu}(\beta), \mathbf{x}_0)$ . We only need to consider  $\beta \in B_{\delta_n}(\beta^*)$ . For  $\beta \in B_{\delta_n}(\beta^*)$ , we have

$$\begin{aligned} & \sqrt{n}(\mu(\beta, \mathbf{x}_0) - \mu(\beta^*, \mathbf{x}_0)) \\ &= \sqrt{n}(\mu(\beta, \mathbf{x}_0) - \mu(\beta_{\gamma^*}, \mathbf{x}_0) + \mu(\tilde{\nu}(\beta_{\gamma^*}), \mathbf{x}_0) - \mu(\beta^*, \mathbf{x}_0)). \end{aligned}$$

Since  $\beta \in B_{\delta_n}(\beta^*)$ , for  $i \notin \gamma^*$ ,  $|\beta_i| < 2\sigma_{0,n} \log(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}})$ ; and for  $i \in \gamma^*$ ,  $|\tilde{\nu}(\beta)_i - \beta_i^*| < \delta \lesssim \frac{1}{\sqrt[3]{nr_n}}$ . Therefore,

$$|\sqrt{n}(\mu(\beta, \mathbf{x}_0) - \mu(\beta_{\gamma^*}, \mathbf{x}_0))| = |\sqrt{n} \sum_{i \notin \gamma^*} \beta_i (\mu_i(\tilde{\beta}, \mathbf{x}_0))| \leq \sqrt{n} K_n M 2\sigma_{0,n} \log(\frac{\sigma_{1,n}}{\lambda_n \sigma_{0,n}}) = o(1),$$

where  $\mu_i(\beta, \mathbf{x}_0)$  denotes the first derivative of  $\mu(\beta, \mathbf{x}_0)$  with respect to the  $i$ th component of  $\beta$ , and  $\tilde{\beta}$  denotes a point between  $\beta$  and  $\beta_{\gamma^*}$ . Further,

$$\begin{aligned} & \mu(\tilde{\nu}(\beta_{\gamma^*}), \mathbf{x}_0) - \mu(\beta^*, \mathbf{x}_0) \\ &= \sqrt{n} \sum_{i \in \gamma^*} (\tilde{\nu}(\beta)_i - \beta_i^*) \mu_i(\beta^*, \mathbf{x}_0) + \sqrt{n} \sum_{i \in \gamma^*} \sum_{j \in \gamma^*} (\tilde{\nu}(\beta)_i - \beta_i^*) \mu_{i,j}(\tilde{\beta}, \mathbf{x}_0) (\tilde{\nu}(\beta)_j - \beta_j^*) \\ &= \sqrt{n} \sum_{i \in \gamma^*} ((\tilde{\nu}(\beta)_i - \beta_i^*) \mu_i(\beta^*, \mathbf{x}_0) + o(1)), \end{aligned}$$

where  $\mu_{i,j}(\beta, \mathbf{x}_0)$  denotes the second derivative of  $\mu(\beta, \mathbf{x}_0)$  with respect to the  $i$ th and  $j$ th components of  $\beta$  and  $\tilde{\beta}$  is a point between  $\tilde{\nu}(\beta)$  and  $\beta^*$ . Summarizing the above two equations, we have

$$\sqrt{n}(\mu(\beta, \mathbf{x}_0) - \mu(\beta^*, \mathbf{x}_0)) = \sqrt{n} \sum_{i \in \gamma^*} ((\tilde{\nu}(\beta)_i - \beta_i^*) \mu_i(\beta^*, \mathbf{x}_0) + o(1)).$$

By Theorem 2.1,  $\pi[\sqrt{n}(\tilde{\nu}(\beta) - \beta^*) \mid D_n] \rightsquigarrow N(0, \mathbf{V})$ , where  $\mathbf{V} = (v_{ij})$ , and  $v_{i,j} = E(h^{i,j}(\beta^*))$  if  $i, j \in \gamma^*$  and 0 otherwise. Then we have  $\pi[\sqrt{n}(\mu(\beta, \mathbf{x}_0) - \mu(\beta^*, \mathbf{x}_0)) \mid D_n] \rightsquigarrow N(0, \Sigma)$ , where  $\Sigma = \nabla_{\gamma^*} \mu(\beta^*, \mathbf{x}_0)^T H^{-1} \nabla_{\gamma^*} \mu(\beta^*, \mathbf{x}_0)$  and  $H = E(-\nabla_{\gamma^*}^2 l_n(\beta^*))$ .

## 3 Theory of Prior Annealing: Proof of Theorem 3.1

Our proof follows the proof of Theorem 2 in Chen et al. (2015). SGLD use the first order integrator (see Lemma 12 of Chen et al. (2015) for the detail). Then we have

$$\begin{aligned} \mathbb{E}(\psi(\beta^{(t+1)})) &= \psi(\beta^{(t)}) + \epsilon_t \mathcal{L}_t \psi(\beta^{(t)}) + O(\epsilon_t^2) \\ &= \psi(\beta^{(t)}) + \epsilon_t (\mathcal{L}_t - \mathcal{L}) \psi(\beta^{(t)}) + \epsilon_t \mathcal{L} \psi(\beta^{(t)}) + O(\epsilon_t^2). \end{aligned}$$

Note that by Poisson equation,  $\mathcal{L} \psi(\beta) = \phi(\beta) - \int \phi(\beta) \pi(\beta \mid D_n, \eta^*, \sigma_{0,n}^*) d\beta$ . Taking expectation on both sides of the equation, summing over  $t = 0, 1, \dots, T-1$ , and dividing  $\epsilon T$  on both sides of

the equation, we have

$$\begin{aligned} & \mathbb{E} \left( \frac{1}{T} \sum_{t=1}^{T-1} \phi(\boldsymbol{\beta}^{(t)}) - \int \phi(\boldsymbol{\beta}) \pi(\boldsymbol{\beta} | D_n, \eta^*, \sigma_{0,n}^*) \right) \\ &= \frac{1}{T\epsilon} (\mathbb{E}(\psi(\boldsymbol{\beta}^{(T)})) - \psi(\boldsymbol{\beta}^{(0)})) - \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}(\delta_t \psi(\boldsymbol{\beta}^{(t)})) + O(\epsilon). \end{aligned}$$

To characterize the order of  $\delta_t = \mathcal{L}_t - \mathcal{L}$ , we first study the difference of the drift term

$$\begin{aligned} & \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | D_{m,n}^{(t)}, \eta^{(t)}, \sigma_{0,n}^{(t)})) - \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | D_n, \eta^*, \sigma_{0,n}^*)) \\ &= \sum_{i=1}^n \nabla \log(p_{\boldsymbol{\beta}^{(t)}}(\mathbf{x}_i, y_i)) - \frac{n}{m} \sum_{j=1}^m \nabla \log(p_{\boldsymbol{\beta}^{(t)}}(\mathbf{x}_{i_j}, y_{i_j})) \\ & \quad + \eta^{(t)} \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^{(t)}, \sigma_{1,n})) - \eta^* \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^*, \sigma_{1,n})). \end{aligned}$$

Use of the mini-batch data gives an unbiased estimator of the full gradient, i.e.

$$\mathbb{E} \left( \sum_{i=1}^n \nabla \log(p_{\boldsymbol{\beta}^{(t)}}(\mathbf{x}_i, y_i)) - \frac{n}{m} \sum_{j=1}^m \nabla \log(p_{\boldsymbol{\beta}^{(t)}}(\mathbf{x}_{i_j}, y_{i_j})) \right) = 0.$$

For the prior part, let  $p(\sigma)$  denote the density function of  $N(0, \sigma)$ . Then we have

$$\begin{aligned} & \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^{(t)}, \sigma_{1,n})) \\ &= - \frac{(1 - \lambda_n) p(\sigma_{0,n}^{(t)})}{(1 - \lambda_n) p(\sigma_{0,n}^{(t)}) + \lambda_n p(\sigma_{1,n})} \frac{\boldsymbol{\beta}^{(t)}}{\sigma_{0,n}^{(t)2}} - \frac{\lambda_n p(\sigma_{1,n})}{(1 - \lambda_n) p(\sigma_{0,n}^{(t)}) + \lambda_n p(\sigma_{1,n})} \frac{\boldsymbol{\beta}^{(t)}}{\sigma_{1,n}^2}, \end{aligned}$$

and thus  $\mathbb{E} |\nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^{(t)}, \sigma_{1,n}))| \leq \frac{2\mathbb{E}|\boldsymbol{\beta}^{(t)}|}{\sigma_{0,n}^*}$ . By Assumption 5.2, we have

$$\begin{aligned} & \mathbb{E} (|\eta^{(t)} \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^{(t)}, \sigma_{1,n})) - \eta^* \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^*, \sigma_{1,n}))|) \\ &= \mathbb{E} (|\eta^{(t)} \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^{(t)}, \sigma_{1,n})) - \eta^* \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^{(t)}, \sigma_{1,n}))|) \\ & \quad + \mathbb{E} (|\eta^* \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^{(t)}, \sigma_{1,n})) - \eta^* \nabla \log(\pi(\boldsymbol{\beta}^{(t)} | \lambda_n, \sigma_{0,n}^*, \sigma_{1,n}))|) \\ &\leq \frac{2M}{\sigma_{0,n}^*} |\eta^{(t)} - \eta^*| + \eta^* M |\sigma_{0,n}^{(t)} - \sigma_{0,n}^*|. \end{aligned}$$

By Assumption 5.1,  $\mathbb{E}(\psi(\boldsymbol{\beta}^{(t)})) \leq \infty$ . Then

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}(\delta_t \psi(\boldsymbol{\beta}^{(t)})) = O \left( \frac{1}{T} \sum_{t=0}^{T-1} (|\eta^{(t)} - \eta^*| + |\sigma_{0,n}^{(t)} - \sigma_{0,n}^*|) \right).$$

Note that by assumption 5.1,  $|\psi(\boldsymbol{\beta}^{(T)}) - \psi(\boldsymbol{\beta}^{(0)})|$  is bounded. Then

$$\mathbb{E} \left( \frac{1}{T} \sum_{t=1}^{T-1} \phi(X_t) - \int \phi(\boldsymbol{\beta}) \pi(\boldsymbol{\beta} | D_n, \eta^*, \sigma_{0,n}^*) \right) = O \left( \frac{1}{T\epsilon} + \frac{\sum_{t=0}^{T-1} (|\eta^{(t)} - \eta^*| + |\sigma_{0,n}^{(t)} - \sigma_{0,n}^*|)}{T} + \epsilon \right).$$

## 4 Construct Confidence Interval

Theorem 2.2 implies that a faithful prediction interval can be constructed for the sparse neural network learned by the proposed algorithms. In practice, for a normal regression problem with noise  $N(0, \sigma^2)$ , to construct the prediction interval for a test point  $\mathbf{x}_0$ , the terms  $\sigma^2$  and  $\Sigma = \nabla_{\boldsymbol{\gamma}^*} \mu(\boldsymbol{\beta}^*, \mathbf{x}_0)^T H^{-1} \nabla_{\boldsymbol{\gamma}^*} \mu(\boldsymbol{\beta}^*, \mathbf{x}_0)$  in Theorem 2.2 need to be estimated from data. Let  $D_n = (\mathbf{x}^{(i)}, y^{(i)})_{i=1, \dots, n}$  be the training set and  $\mu(\boldsymbol{\beta}, \cdot)$  be the predictor of the network model with parameter  $\boldsymbol{\beta}$ . We can follow the following procedure to construct the prediction interval for the test point  $\mathbf{x}_0$ :

- Run algorithm 1, let  $\hat{\beta}$  be an estimation of the network parameter at the end of the algorithm and  $\hat{\gamma}$  be the corresponding network structure.
- Estimate  $\sigma^2$  by

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\mu(\hat{\beta}, \mathbf{x}^{(i)}) - y^{(i)})^2.$$

- Estimate  $\Sigma$  by

$$\hat{\Sigma} = \nabla_{\hat{\gamma}} \mu(\hat{\beta}, \mathbf{x}_0)^T (-\nabla_{\hat{\gamma}}^2 l_n(\hat{\beta}))^{-1} \nabla_{\hat{\gamma}} \mu(\hat{\beta}, \mathbf{x}_0).$$

- Construct the prediction interval as

$$\left( \mu(\hat{\beta}, \mathbf{x}_0) - 1.96 \sqrt{\frac{1}{n} \hat{\Sigma} + \hat{\sigma}^2}, \mu(\hat{\beta}, \mathbf{x}_0) + 1.96 \sqrt{\frac{1}{n} \hat{\Sigma} + \hat{\sigma}^2} \right).$$

Here, by the structure selection consistency (Lemma 2.2) and consistency of the MLE for the learnt structure Portnoy (1988), we replace  $\beta^*$  and  $\gamma^*$  in Theorem 2.2 by  $\hat{\beta}$  and  $\hat{\gamma}$ .

If the dimension of the sparse network is still too high and the computation of  $\hat{\Sigma}$  becomes prohibitive, the following Bayesian approach can be used to construct confidence intervals.

- Running SGMCMC algorithm to get a sequence of posterior samples:  $\beta^{(1)}, \dots, \beta^{(m)}$ .
- Estimating  $\sigma^2$  by  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \mu^{(i)})^2$ , where

$$\mu^{(i)} = \frac{1}{m} \sum_{j=1}^m \mu(\beta^{(j)}, \mathbf{x}^{(i)}), i = 1, \dots, n,$$

- Estimate the prediction mean by

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m \mu(\beta^{(i)}, \mathbf{x}_0).$$

- Estimate the prediction variance by

$$\hat{V} = \frac{1}{m} \sum_{i=1}^m (\mu(\beta^{(i)}, \mathbf{x}_0) - \hat{\mu})^2 + \hat{\sigma}^2.$$

- Construct the prediction interval as

$$(\mu - 1.96\sqrt{\hat{V}}, \mu + 1.96\sqrt{\hat{V}}).$$

## 5 Prior Annealing

In this section, we give some graphical illustration of the prior annealing algorithm. In practice, the negative log-prior puts penalty on parameter weights. The mixture Gaussian prior behaves like a piecewise  $L_2$  penalty with different weights on different regions. Figure 1 shows the shape of a negative log-mixture Gaussian prior. In step (iii) of Algorithm 1, the condition  $\pi(\gamma_i = 1|\beta_i) > 0.5$  splits the parameters into two parts. For the  $\beta_i$ 's with large magnitudes, the slab component  $N(0, \sigma_{1,n}^2)$  plays the major role in the prior, imposing a small penalty on the parameter. For the  $\beta_i$ 's with smaller magnitudes, the spike component  $N(0, \sigma_{0,n}^2)$  plays the major role in the prior, imposing a large penalty on the parameters to push them toward zero in training.

Figure 2 shows the shape of negative log-prior and  $\pi(\gamma_i = 1|\beta_i)$  for different choices of  $\sigma_{0,n}^2$  and  $\lambda_n$ . As we can see from the plot,  $\sigma_{0,n}^2$  plays the major role in determining the effect of the prior. Let  $\alpha$  be the threshold in step (iii) of Algorithm 1 of the main body, i.e. the positive solution to  $\pi(\gamma_i = 1|\beta_i) = 0.5$ . In general, a smaller  $\sigma_{0,n}^2$  will result in a smaller  $\alpha$ . If a very small  $\sigma_{0,n}^2$  is used

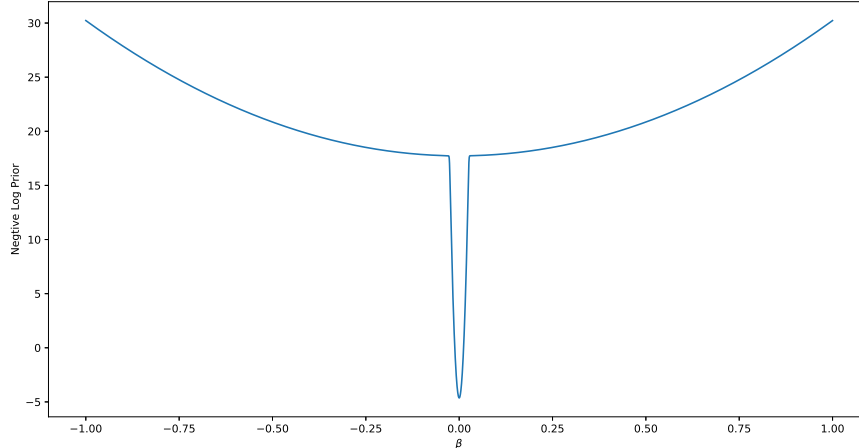


Figure 1: Negative logarithm of the mixture Gaussian prior.

in the prior from the beginning, then most of  $\beta_i$ 's at initialization will have a magnitude larger than  $\alpha$  and the slab component  $N(0, \sigma_{1,n}^2)$  of the prior will dominate most parameters. As a result, it will be difficult to find the desired sparse structure. Following the proposed prior annealing procedure, we can start with a larger  $\sigma_{0,n}^2$ , i.e. a larger threshold  $\alpha$  and a relatively smaller penalty for those  $|\beta_i| < \alpha$ . As we gradually decrease the value of  $\sigma_{0,n}^2$ ,  $\alpha$  decreases, and the penalty imposed on the small weights increases and drives them toward zero. The prior annealing allows us to gradually sparsify the DNN and impose more and more penalties on the parameters close to 0.

## 6 Experimental Setups

### 6.1 Simulated examples

**Prior annealing** We follow simple implementation of Algorithm given in section 3.1. We run SGHMC for  $T = 80000$  iterations with constant learning rate  $\epsilon_t = 0.001$ , momentum  $1 - \alpha = 0.9$  and subsample size  $m = 500$ . We set  $\lambda_n = 1e - 7$ ,  $\sigma_{1,n}^2 = 1e - 2$ ,  $(\sigma_{0,n}^{init})^2 = 5e - 5$ ,  $(\sigma_{0,n}^{end})^2 = 1e - 6$  and  $T_1 = 5000$ ,  $T_2 = 20000$ ,  $T_3 = 60000$ . We set temperature  $\tau = 0.1$  for  $t < T_3$  and for  $t > T_3$ , we gradually decrease temperature  $\tau$  by  $\tau = \frac{0.1}{t - T_3}$ . After structure selection, the model is fine tuned for 40000 iterations. The number of iteration setup is the same as Sun et al. (2021).

**Other Methods** Spinn, Dropout and DNN are trained with the same network structure using SGD with momentum. Same as our method, we use constant learning rate 0.001, momentum 0.9, subsample size 500 and traing the model for 80000 iterations. For Spinn, we use LASSO penalty and the regularization parameter is selected from  $\{0.05, 0.06, \dots, 0.15\}$  according to the performance on validation data set. For Dropout, the dropout rate is set to be 0.2 for the first layer and 0.5 for the other layers. Other baseline methods BART50, LASSO, SIS are implemented using R-package *randomForest*, *glmnet*, *BART* and *SIS* respectively with default parameters.

### 6.2 CIFAR10

We follow the standard training procedure as in Lin et al. (2020), i.e. we train the model with SGHMC for  $T = 300$  epochs, with initial learning rate  $\epsilon_0 = 0.1$ , momentum  $1 - \alpha = 0.9$ , temperature  $\tau = 0.001$ , mini-batch size  $m = 128$ . The learning rate is divided by 10 at 150th and 225th epoch. We follow the implementation given in section 3.1 and use  $T_1 = 150$ ,  $T_2 = 200$ ,  $T_3 = 225$ , where  $T_i$ s are number of epochs. We set temperature  $\tau = 0.01$  for  $t < T_3$  and gradually decrease  $\tau$  by  $\tau = \frac{0.01}{t - T_3}$  for  $t > T_3$ . We set  $\sigma_{1,n}^2 = 0.04$  and  $(\sigma_{0,n}^{init})^2 = 10 \times (\sigma_{0,n}^{end})^2$  and use different  $\sigma_{0,n}^{end}$ ,  $\lambda_n$  for different network size and target sparsity level. The detailed settings are given below:



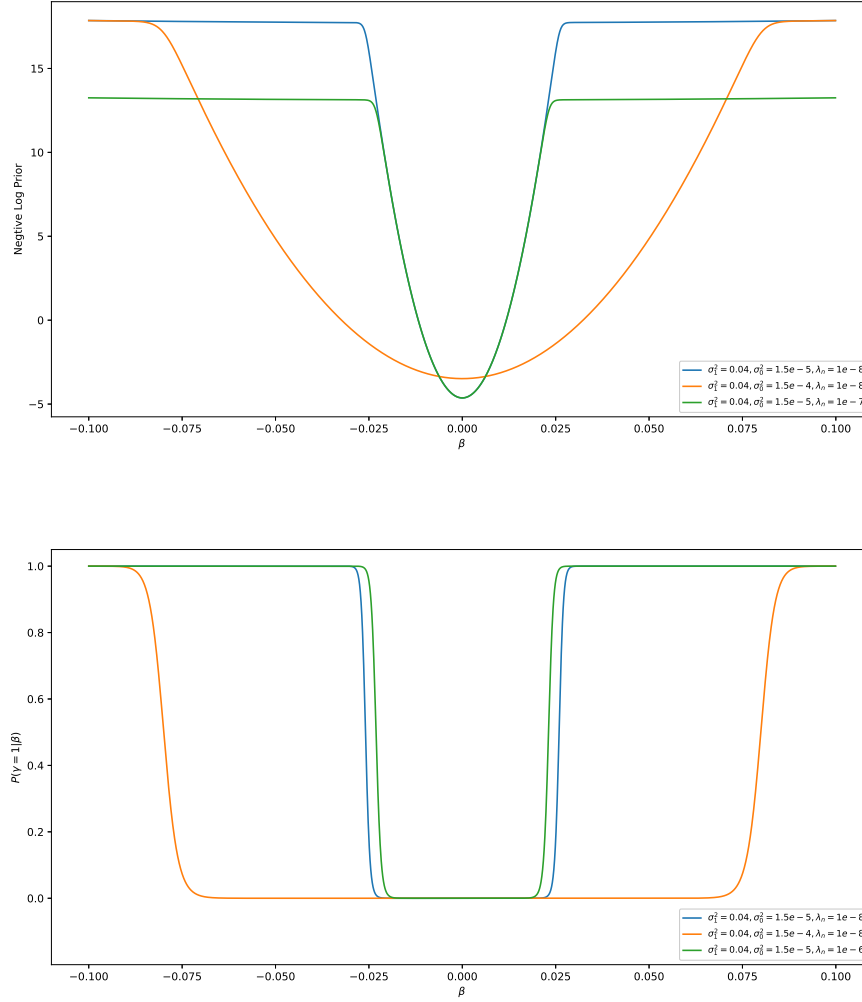


Figure 2: Negative log-prior and  $\pi(\gamma = 1|\beta)$  for different choices of  $\sigma_{0,n}^2$  and  $\lambda_n$ .

- ResNet20 with target sparsity level 20%:  $(\sigma_{0,n}^{end})^2 = 1.5e - 5, \lambda_n = 1e - 8$
- ResNet20 with target sparsity level 10%:  $(\sigma_{0,n}^{end})^2 = 6e - 5, \lambda_n = 1e - 9$
- ResNet32 with target sparsity level 10%:  $(\sigma_{0,n}^{end})^2 = 3e - 5, \lambda_n = 2e - 9$
- ResNet32 with target sparsity level 5%:  $(\sigma_{0,n}^{end})^2 = 1e - 4, \lambda_n = 2e - 8$

## References

- Castillo, I. and Rousseau, J. (2015). Supplement to “a bernstein–von mises theorem for smooth functionals in semiparametric models”. *Annals of Statistics*, 43(6):2353–2383.
- Chen, C., Ding, N., and Carin, L. (2015). On the convergence of stochastic gradient mcmc algorithms with high-order integrators. In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, pages 2278–2286.
- Fefferman, C. (1994). Reconstructing a neural net from its output. *Revista Matemática Iberoamericana*, 10(3):507–555.

- Feng, J. and Simon, N. (2017). Sparse-input neural networks for high-dimensional nonparametric regression and classification. *arXiv preprint arXiv:1711.07592*.
- Lin, T., Stich, S. U., Barba, L., Dmitriev, D., and Jaggi, M. (2020). Dynamic model pruning with feedback. In *International Conference on Learning Representations*.
- Portnoy, S. (1988). Asymptotic behavior of likelihood methods for exponential families when the number of parameters tend to infinity. *The Annals of Statistics*, 16(1):356–366.
- Sun, Y., Song, Q., and Liang, F. (2021). Consistent sparse deep learning: Theory and computation. *Journal of the American Statistical Association*, page in press.