

---

# Supplimentary Material for *An Adaptive Empirical Bayesian Method for Sparse Deep Learning*

---

**Wei Deng**  
Purdue University  
deng106@purdue.edu

**Xiao Zhang**  
Purdue University  
zhang923@purdue.edu

**Faming Liang**  
Purdue University  
fmliang@purdue.edu

**Guang Lin**  
Purdue University  
guanglin@purdue.edu

In this supplementary material, we review the related methodologies in §1, prove the convergence in §2, present additional simulation of logistic regression in §3, illustrate more regression examples on UCI datasets in §4, and show the experimental setup in §5.

## 1 Stochastic Approximation

### 1.1 Special Case: Robbins–Monro Algorithm

Robbins–Monro algorithm is the first stochastic approximation algorithm to deal with the root finding problem which also applies to the stochastic optimization problem. Given the random output of  $H(\theta, \beta)$  with respect to  $\beta$ , our goal is to find  $\theta^*$  such that

$$h(\theta^*) = \mathbb{E}_{\theta^*}[H(\theta^*, \beta)] = \int H(\theta^*, \beta) f_{\theta^*}(d\beta) = 0, \quad (1)$$

where  $\mathbb{E}_{\theta^*}$  denotes the expectation with respect to the distribution of  $\beta$  given  $\theta^*$ . To implement the Robbins–Monro Algorithm, we can generate iterates as follows\*:

- (1) Sample  $\beta_{k+1}$  from the invariant distribution  $f_{\theta_k}(\beta)$ ,
- (2) Update  $\theta_{k+1} = \theta_k + \omega_{k+1} H(\theta_k, \beta_{k+1})$ .

Note that in this algorithm,  $H(\theta, \beta)$  is the unbiased estimator of  $h(\theta)$ , that is for  $k \in \mathbb{N}^+$ , we have

$$\mathbb{E}_{\theta_k}[H(\theta_k, \beta_{k+1}) - h(\theta_k) | \mathcal{F}_k] = 0. \quad (2)$$

If there exists an antiderivative  $Q(\theta, \beta)$  that satisfies  $H(\theta, \beta) = \nabla_{\theta} Q(\theta, \beta)$  and  $E_{\theta}[Q(\theta, \beta)]$  is concave, it is equivalent to solving the stochastic optimization problem  $\max_{\theta \in \Theta} E_{\theta}[Q(\theta, \beta)]$ .

### 1.2 General Stochastic Approximation

The stochastic approximation algorithm is an iterative recursive algorithm consisting of two steps:

- (1) Sample  $\beta_{k+1}$  from the transition kernel  $\Pi_{\theta_k}(\beta_k, \cdot)$ , which admits  $f_{\theta_k}(\beta)$  as the invariant distribution,
- (2) Update  $\theta_{k+1} = \theta_k + \omega_{k+1} H(\theta_k, \beta_{k+1})$ .

The general stochastic approximation [Benveniste et al., 1990] differs from the Robbins-Monro algorithm in that sampling  $x$  from a transition kernel instead of a distribution introduces a Markov state-dependent noise  $H(\theta_k, x_{k+1}) - h(\theta_k)$ .

---

\*We change the notation a little bit, where  $\beta_k \in \mathbb{R}^d$  and  $\theta_k$  are the parameters at the  $k$ -th iteration.

## 2 Convergence Analysis

### 2.1 Convergence of Hidden Variables

The stochastic gradient Langevin Dynamics with a stochastic approximation adaptation (SGLD-SA) is a mixed half-optimization-half-sampling algorithm to handle complex Bayesian posterior with latent variables, e.g. the conjugate spike-slab hierarchical prior formulation. Each iteration of the algorithm consists of the following steps:

- (1) Sample  $\beta_{k+1}$  using SGLD based on  $\theta_k$ , i.e.

$$\beta_{k+1} = \beta_k + \epsilon \nabla_{\beta} \tilde{L}(\beta_k, \theta_k) + \sqrt{2\epsilon\tau^{-1}} \eta_k, \quad (3)$$

where  $\eta_k \sim \mathcal{N}(0, I)$ ;

- (2) Optimize  $\theta_{k+1}$  from the following recursion

$$\begin{aligned} \theta_{k+1} &= \theta_k + \omega_{k+1} (g_{\theta_k}(\beta_{k+1}) - \theta_k) \\ &= (1 - \omega_{k+1}) \theta_k + \omega_{k+1} g_{\theta_k}(\beta_{k+1}), \end{aligned} \quad (4)$$

where  $g_{\theta_k}(\cdot)$  is some mapping to derive the optimal  $\theta$  based on the current  $\beta$ .

**Remark:** Define  $H(\theta_k, \beta_{k+1}) = g_{\theta_k}(\beta_{k+1}) - \theta_k$ . In this formulation, our target is to find  $\theta^*$  that solves  $h(\theta^*) = \mathbb{E}[H(\theta, \beta)] = 0$ .

#### General Assumptions

To provide the  $L_2$  upper bound for SGLD-SA, we first lay out the following assumptions:

**Assumption 1** (Step size and Convexity).  $\{\omega_k\}_{k \in \mathbb{N}}$  is a positive decreasing sequence of real numbers such that

$$\omega_k \rightarrow 0, \quad \sum_{k=1}^{\infty} \omega_k = +\infty. \quad (5)$$

There exist  $\delta > 0$  and  $\theta^*$  such that for  $\theta \in \Theta$ :<sup>†</sup>

$$\langle \theta - \theta^*, h(\theta) \rangle \leq -\delta \|\theta - \theta^*\|^2, \quad (6)$$

with additionally

$$\lim_{k \rightarrow \infty} \inf 2\delta \frac{\omega_k}{\omega_{k+1}} + \frac{\omega_{k+1} - \omega_k}{\omega_{k+1}^2} > 0. \quad (7)$$

Then for any  $\alpha \in (0, 1]$  and suitable  $A$  and  $B$ , a practical  $\omega_k$  can be set as

$$\omega_k = A(k + B)^{-\alpha} \quad (8)$$

**Assumption 2** (Smoothness).  $L(\beta, \theta)$  is  $M$ -smooth with  $M > 0$ , i.e. for any  $\beta, \iota \in \mathcal{B}$ ,  $\theta, \nu \in \Theta$ .

$$\|\nabla_{\beta} L(\beta, \theta) - \nabla_{\beta} L(\iota, \nu)\| \leq M \|\beta - \iota\| + M \|\theta - \nu\|. \quad (9)$$

**Assumption 3** (Dissipative). There exist constants  $m > 0, b \geq 0$ , s.t. for all  $\beta \in \mathcal{B}$  and  $\theta \in \Theta$ , we have

$$\langle \nabla_{\beta} L(\beta, \theta), \beta \rangle \leq b - m \|\beta\|^2. \quad (10)$$

**Assumption 4** (Gradient condition). The stochastic noise  $\chi_k \in \mathcal{B}$ , which comes from  $\nabla_{\beta} \tilde{L}(\beta_k, \theta_k) - \nabla_{\beta} L(\beta_k, \theta_k)$ , is a white noise or Martingale difference noise and is independent with each other.

$$\mathbb{E}[\chi_k | \mathcal{F}_k] = 0. \quad (11)$$

The scale of the noise is bounded by

$$\mathbb{E} \|\chi\|^2 \leq M^2 \mathbb{E} \|\beta\|^2 + M^2 \mathbb{E} \|\theta\|^2 + B^2. \quad (12)$$

for constants  $M, B > 0$ .

---

<sup>†</sup>  $\|\cdot\|$  is short for  $\|\cdot\|_2$

In addition to the assumptions, we also assume the existence of Markov transition kernel, the proof goes beyond the scope of our paper.

**Proposition 1.** *There exist constants  $M, B > 0$  such that*

$$\|g_{\theta}(\beta)\|^2 \leq M^2\|\beta\|^2 + B^2 \quad (13)$$

*Proof.* As shown in Eq.(12), Eq.(13) and Eq.(15) in the main body,  $\rho, \delta$  and  $\kappa$  are clearly bounded. It is also easy to verify that  $\sigma$  in Eq.(14) in the main body satisfies (13). For convenience, we choose the same  $M$  and  $B$  (large enough) as in (12).  $\square$

**Proposition 2.** *For any  $\beta \in B$ , it holds that*

$$\|\nabla_{\beta}L(\beta, \theta)\|^2 \leq 3M^2\|\beta\|^2 + 3M^2\|\theta\|^2 + 3B^2 \quad (14)$$

for constants  $M$  and  $B$ .

*Proof.* Suppose there is a minimizer  $(\theta^*, \beta^*)$  such that  $\nabla_{\beta}L(\beta^*, \theta^*) = 0$  and  $\theta^*$  has reached the stationary point, following Assumption 3 we have,

$$\langle \nabla_{\beta}L(\beta^*, \theta^*), \beta^* \rangle \leq b - m\|\beta^*\|^2.$$

Therefore,  $\|\beta^*\|^2 \leq \frac{b}{m}$ . Since  $\theta^*$  is the stationary point,  $\theta^* = (1 - \omega)\theta^* + \omega g_{\theta^*}(\beta^*)$ . By (13), we have  $\|g_{\theta^*}(\beta^*)\|^2 \leq M^2\|\beta^*\|^2 + B^2$ , which implies that  $\|\theta^*\|^2 = \|g_{\theta^*}(\beta^*)\|^2 \leq M^2\|\beta^*\|^2 + B^2 \leq \frac{b}{m}M^2 + B^2$ . By the smoothness assumption 2, we have

$$\begin{aligned} & \|\nabla_{\beta}L(\beta, \theta)\| \\ & \leq \|\nabla_{\beta}L(\beta^*, \theta^*)\| + M\|\beta - \beta^*\| + M\|\theta - \theta^*\| \\ & \leq 0 + M(\|\beta\| + \sqrt{\frac{b}{m}} + \|\theta\| + \|\theta^*\|) \\ & \leq M\|\theta\| + M\|\beta\| + M(\sqrt{\frac{b}{m}} + \sqrt{\frac{b}{m}M^2 + B^2}) \\ & \leq M\|\theta\| + M\|\beta\| + \bar{B}, \end{aligned}$$

where  $\bar{B} = M(\sqrt{\frac{b}{m}M^2 + B^2} + \sqrt{\frac{b}{m}})$ . Therefore,

$$\|\nabla_{\beta}L(\beta, \theta)\|^2 \leq 3M^2\|\beta\|^2 + 3M^2\|\theta\|^2 + 3\bar{B}^2.$$

For notation simplicity, we can choose the same  $B$  (large enough) to bound (12), (13) and (14).  $\square$

**Lemma 1** (Uniform  $L_2$  bounds). *For all  $0 < \epsilon < \text{Re}(\frac{m - \sqrt{m^2 - 4M^2(M^2 + 1)}}{4M^2(M^2 + 1)})$ , there exist  $G, \bar{G} > 0$  such that  $\sup \mathbb{E}\|\beta_k\|^2 \leq G$  and  $\sup \mathbb{E}\|\theta_k\|^2 \leq \bar{G}$ , where  $G = \|\beta_0\|^2 + \frac{1}{m}(b + 2\epsilon B^2(M^2 + 1) + \tau d)$  and  $\bar{G} = M^2G + B^2$ .*

*Proof.* From (3), we have

$$\begin{aligned} & \mathbb{E}\|\beta_{k+1}\|^2 \\ & = \mathbb{E}\left\|\beta_k + \epsilon \nabla_{\beta} \tilde{L}(\beta_k, \theta_k)\right\|^2 + 2\tau\epsilon \mathbb{E}\|\eta_k\|^2 + \sqrt{8\epsilon\tau} \mathbb{E}\langle \beta_k + \epsilon \nabla_{\beta} \tilde{L}(\beta_k, \theta_k), \eta_k \rangle \\ & = \mathbb{E}\left\|\beta_k + \epsilon \nabla_{\beta} \tilde{L}(\beta_k, \theta_k)\right\|^2 + 2\tau\epsilon d, \end{aligned} \quad (15)$$

Moreover, the first item in (15) can be expanded to

$$\begin{aligned} & \mathbb{E}\left\|\beta_k + \epsilon \nabla_{\beta} \tilde{L}(\beta_k, \theta_k)\right\|^2 \\ & = \mathbb{E}\|\beta_k + \epsilon \nabla_{\beta}L(\beta_k, \theta_k)\|^2 + \epsilon^2 \mathbb{E}\|\chi_k\|^2 - 2\epsilon \mathbb{E}[\mathbb{E}(\langle \beta_k + \epsilon \nabla_{\beta}L(\beta_k, \theta_k), \chi_k \rangle | \mathcal{F}_k)] \\ & = \mathbb{E}\|\beta_k + \epsilon \nabla_{\beta}L(\beta_k, \theta_k)\|^2 + \epsilon^2 \mathbb{E}\|\chi_k\|^2, \end{aligned} \quad (16)$$

where (11) is used to cancel the inner product item.

Turning to the first item of (16), the dissipativity condition (10) and the boundness of  $\nabla_{\beta}L(\beta, \theta)$  (14) give us:

$$\begin{aligned}
& \mathbb{E} \|\beta_k + \epsilon \nabla_{\beta}L(\beta_k, \theta_k)\|^2 \\
&= \mathbb{E} \|\beta_k\|^2 + 2\epsilon \mathbb{E} \langle \beta_k, \nabla_{\beta}L(\beta_k, \theta_k) \rangle + \epsilon^2 \mathbb{E} \|\nabla_{\beta}L(\beta_k, \theta_k)\|^2 \\
&\leq \mathbb{E} \|\beta_k\|^2 + 2\epsilon(b - m\mathbb{E} \|\beta_k\|^2) + \epsilon^2(3M^2\mathbb{E} \|\beta_k\|^2 + 3M^2\mathbb{E} \|\theta_k\|^2 + 3B^2) \\
&= (1 - 2\epsilon m + 3\epsilon^2 M^2) \mathbb{E} \|\beta_k\|^2 + 2\epsilon b + 3\epsilon^2 B^2 + 3\epsilon^2 M^2 \mathbb{E} \|\theta_k\|^2.
\end{aligned} \tag{17}$$

By (12), the second item of (16) is bounded by

$$\mathbb{E} \|\chi_k\|^2 \leq M^2 \mathbb{E} \|\beta_k\|^2 + M^2 \mathbb{E} \|\theta_k\|^2 + B^2. \tag{18}$$

Combining (15), (16), (17) and (18), we have

$$\mathbb{E} \|\beta_{k+1}\|^2 \leq (1 - 2\epsilon m + 4\epsilon^2 M^2) \mathbb{E} \|\beta_k\|^2 + 2\epsilon b + 4\epsilon^2 B^2 + 4\epsilon^2 M^2 \mathbb{E} \|\theta_k\|^2 + 2\tau\epsilon d. \tag{19}$$

Next we use proof by induction to show for  $k = 1, 2, \dots, \infty$ ,  $\mathbb{E} \|\beta_k\|^2 \leq G$ , where

$$G = \mathbb{E} \|\beta_0\|^2 + \frac{b + 2\epsilon B^2(M^2 + 1) + \tau d}{m - 2\epsilon M^2(M^2 + 1)}. \tag{20}$$

First of all, the case of  $k = 0, 1$  is trivial. Then if we assume for each  $k \in 2, 3, \dots, t$ ,  $\mathbb{E} \|\beta_k\|^2 \leq G$ ,  $\mathbb{E} \|g(\beta_k)\|^2 \leq M^2 G + B^2$ ,  $\mathbb{E} \|\theta_{k-1}\|^2 \leq M^2 G + B^2$ . It follows that,

$$\begin{aligned}
& \mathbb{E} \|\theta_k\|^2 = \mathbb{E} \|(1 - \omega_k)\theta_{k-1} + \omega_k g(\beta_k)\|^2 \\
&\leq (1 - \omega_k)^2 \mathbb{E} \|\theta_{k-1}\|^2 + \omega_k^2 \mathbb{E} \|g(\beta_k)\|^2 + 2(1 - \omega_k)\omega_k \mathbb{E} \langle \theta_{k-1}, g(\beta_k) \rangle \\
&\leq (1 - \omega_k)^2 \mathbb{E} \|\theta_{k-1}\|^2 + \omega_k^2 \mathbb{E} \|g(\beta_k)\|^2 + 2(1 - \omega_k)\omega_k \sqrt{\mathbb{E} \|\theta_{k-1}\|^2 \mathbb{E} \|g(\beta_k)\|^2} \\
&\leq (1 - \omega_k)^2 (M^2 G + B^2) + \omega_k^2 (M^2 G + B^2) + 2(1 - \omega_k)\omega_k (M^2 G + B^2) \\
&= M^2 G + B^2,
\end{aligned}$$

Next, we proceed to prove  $\mathbb{E} \|\beta_{t+1}\|^2 \leq G$  and  $\mathbb{E} \|\theta_{t+1}\|^2 \leq M^2 G + B^2$ . Following (19), we have

$$\begin{aligned}
& \mathbb{E} \|\beta_{t+1}\|^2 \\
&\leq (1 - 2\epsilon m + 4\epsilon^2 M^2) \mathbb{E} \|\beta_t\|^2 + 2\epsilon b + 4\epsilon^2 B^2 + 4\epsilon^2 M^2 \mathbb{E} \|\theta_t\|^2 + 2\tau\epsilon d \\
&\leq (1 - 2\epsilon m + 4\epsilon^2 M^2) G + 2\epsilon b + 4\epsilon^2 B^2 + 4\epsilon^2 M^2 (M^2 G + B^2) + 2\tau\epsilon d \\
&\leq (1 - 2\epsilon m + 4\epsilon^2 M^2 (M^2 + 1)) G + 2\epsilon b + 4\epsilon^2 B^2 (M^2 + 1) + 2\tau\epsilon d
\end{aligned} \tag{21}$$

Consider the quadratic equation  $1 - 2mx + 4M^2(M^2 + 1)x^2 = 0$ . If  $m^2 - 4M^2(M^2 + 1) \geq 0$ , then the smaller root is  $\frac{m - \sqrt{m^2 - 4M^2(M^2 + 1)}}{4M^2(M^2 + 1)}$  which is positive; otherwise the quadratic equation has no real solutions and is always positive. Fix  $\epsilon \in \left(0, \text{Re} \left( \frac{m - \sqrt{m^2 - 4M^2(M^2 + 1)}}{4M^2(M^2 + 1)} \right)\right)$  so that

$$0 < 1 - 2\epsilon m + 4\epsilon^2 M^2 (M^2 + 1) < 1. \tag{22}$$

With (20), we can further bound (21) as follows:

$$\begin{aligned}
& \mathbb{E} \|\beta_{t+1}\|^2 \\
&\leq (1 - 2\epsilon m + 4\epsilon^2 M^2 (M^2 + 1)) (\mathbb{E} \|\beta_0\|^2 + \mathbb{I}) + 2\epsilon b + 4\epsilon^2 B^2 (M^2 + 1) + 2d\tau\epsilon \\
&= (1 - 2\epsilon m + 4\epsilon^2 M^2 (M^2 + 1)) \mathbb{E} \|\beta_0\|^2 + \mathbb{I} - (2\epsilon b + 4\epsilon^2 B^2 (M^2 + 1) + 2d\tau\epsilon) \\
&\quad + (2\epsilon b + 4\epsilon^2 B^2 (M^2 + 1) + 2\epsilon\tau d) \\
&\leq \mathbb{E} \|\beta_0\|^2 + \mathbb{I} \equiv G,
\end{aligned} \tag{23}$$

where  $\mathbb{I} = \frac{b + 2\epsilon B^2(M^2 + 1) + d\tau}{m - 2\epsilon M^2(M^2 + 1)}$ , the second to the last inequality comes from (22).

Moreover, from (13), we also have

$$\begin{aligned} \mathbb{E}\|g(\beta_{t+1})\|^2 &\leq M^2\mathbb{E}\|\beta_{t+1}\|^2 + B^2 \leq M^2G + B^2, \\ \mathbb{E}\|\theta_{t+1}\|^2 &= \mathbb{E}\|(1 - \omega_{t+1})\theta_t + \omega_{t+1}g(\beta_{t+1})\|^2 \\ &\leq (1 - \omega_{t+1})^2\mathbb{E}\|\theta_t\|^2 + \omega_{t+1}^2\mathbb{E}\|g(\beta_{t+1})\|^2 + 2(1 - \omega_{t+1})\omega_{t+1}\mathbb{E}\langle\theta_t, g(\beta_{t+1})\rangle \\ &\leq (1 - \omega_{t+1})^2\mathbb{E}\|\theta_t\|^2 + \omega_{t+1}^2\mathbb{E}\|g(\beta_{t+1})\|^2 + 2(1 - \omega_{t+1})\omega_{t+1}\sqrt{\mathbb{E}\|\theta_t\|^2\mathbb{E}\|g(\beta_{t+1})\|^2} \\ &\leq (1 - \omega_{t+1})^2(M^2G + B^2) + \omega_{t+1}^2(M^2G + B^2) + 2(1 - \omega_{t+1})\omega_{t+1}(M^2G + B^2) \\ &= M^2G + B^2, \end{aligned}$$

Therefore, we have proved that for any  $k \in 1, 2, \dots, \infty$ ,  $\mathbb{E}\|\beta_k\|^2$ ,  $\mathbb{E}\|g(\beta_k)\|^2$  and  $\mathbb{E}\|\theta_k\|^2$  are bounded. Furthermore, we notice that  $G$  can be unified to a constant  $G = \mathbb{E}\|\beta_0\|^2 + \frac{1}{m}(b + 2\epsilon B^2(M^2 + 1) + \tau d)$ .  $\square$

**Assumption 5** (Solution of Poisson equation). *For all  $\theta \in \Theta$ , there exists a function  $\mu_\theta$  on  $\beta$  that solves the Poisson equation  $\mu_\theta(\beta) - \Pi_\theta \mu_\theta(\beta) = H(\theta, \beta) - h(\theta)$ , which follows that*

$$H(\theta_k, \beta_{k+1}) = h(\theta_k) + \mu_{\theta_k}(\beta_{k+1}) - \Pi_{\theta_k} \mu_{\theta_k}(\beta_{k+1}). \quad (24)$$

There exists a constant  $C$  such that for all  $\theta \in \Theta$ ,  $\Pi_\theta \mu$  is bounded, i.e.

$$\|\Pi_\theta \mu_\theta\| \leq C \quad (25)$$

We leave the relaxation of the above assumption for future work.

**Proposition 3.** *There exists a constant  $C_1$  so that*

$$\mathbb{E}_\theta[\|H(\theta, \beta)\|^2] \leq C_1(1 + \|\theta - \theta^*\|^2) \quad (26)$$

*Proof.* By (13), we have

$$\mathbb{E}\|g_\theta(\beta) - \theta\|^2 \leq 2\mathbb{E}\|g_\theta(\beta)\|^2 + 2\|\theta\|^2 \leq 2(M^2\mathbb{E}\|\beta\|^2 + B^2) + 2\|\theta\|^2$$

Since we have proved the  $L_2$  boundness of  $\mathbb{E}\|\beta\|^2$ , choose  $C' = \max(2, 2(M^2\mathbb{E}\|\beta\|^2 + B^2))$ , we have

$$\mathbb{E}_\theta[\|H(\theta, \beta)\|^2] \leq C'(1 + \|\theta\|^2) = C'(1 + \|\theta - \theta^* + \theta^*\|^2) \leq C_1(1 + \|\theta - \theta^*\|^2)$$

$\square$

Lemma 2 is a restatement of Lemma 25 (page 247) from Benveniste et al. [1990].

**Lemma 2.** *Suppose  $k_0$  is an integer which satisfies with*

$$\inf_{k \geq k_0} \frac{\omega_{k+1} - \omega_k}{\omega_k \omega_{k+1}} + 2\delta - \omega_{k+1}C_1 > 0.$$

Then for any  $k > k_0$ , the sequence  $\{\Lambda_k^K\}_{k=k_0, \dots, K}$  defined below is increasing

$$\begin{cases} 2\omega_k \prod_{j=k}^{K-1} (1 - 2\omega_{j+1}\delta + \omega_{j+1}^2 C_1) & \text{if } k < K, \\ 2\omega_k & \text{if } k = K. \end{cases} \quad (27)$$

**Lemma 3.** *There exist  $\lambda_0$  and  $k_0$  such that for all  $\lambda \geq \lambda_0$  and  $k \geq k_0$ , the sequence  $u_k = \lambda\omega_k$  satisfies*

$$u_{k+1} \geq (1 - 2\omega_{k+1}\delta + \omega_{k+1}^2 C_1)u_k + \omega_{k+1}^2 C_1 + \omega_{k+1} \bar{C}_1. \quad (28)$$

*Proof.* Replace  $u_k = \lambda\omega_k$  in (28), we have

$$\lambda\omega_{k+1} \geq (1 - 2\omega_{k+1}\delta + \omega_{k+1}^2 C_1)\lambda\omega_k + \omega_{k+1}^2 C_1 + \omega_{k+1}\bar{C}_1. \quad (29)$$

According to (7) in assumption 1, we denote  $\lim_{k \rightarrow \infty} \inf 2\delta\omega_{k+1}\omega_k + \omega_{k+1} - \omega_k$  by  $\Delta_+$ . Then the above inequality (29) can be simplified as

$$\lambda(\Delta_+ - \omega_{k+1}^2\omega_k C_1) \geq \omega_{k+1}^2 C_1 + \omega_{k+1}\bar{C}_1. \quad (30)$$

Since the LHS increases to  $\Delta_+$  and the RHS decreases to 0 as  $k \rightarrow \infty$ . There exist  $\lambda_0$  and  $k_0$  such that for all  $\lambda > \lambda_0$  and  $k > k_0$ , (30) holds.  $\square$

**Theorem 1** ( $L_2$  convergence rate). *Suppose that Assumptions 1-5 hold, there exists a constant  $\lambda$  such that*

$$\mathbb{E} [\|\boldsymbol{\theta}_k - \boldsymbol{\theta}^*\|^2] \leq \lambda\omega_k,$$

*Proof.* Denote  $\mathbf{T}_k = \boldsymbol{\theta}_k - \boldsymbol{\theta}^*$ , with the help of (4) and Poisson equation (24), we deduce that

$$\begin{aligned} & \|\mathbf{T}_{k+1}\|^2 \\ &= \|\mathbf{T}_k\|^2 + \omega_{k+1}^2 \|H(\boldsymbol{\theta}_k, \boldsymbol{\beta}_{k+1})\|^2 + 2\omega_{k+1} \langle \mathbf{T}_k, H(\boldsymbol{\theta}_k, \boldsymbol{\beta}_{k+1}) \rangle \\ &= \|\mathbf{T}_k\|^2 + \omega_{k+1}^2 \|H(\boldsymbol{\theta}_k, \boldsymbol{\beta}_{k+1})\|^2 + 2\omega_{k+1} \langle \mathbf{T}_k, h(\boldsymbol{\theta}_k) \rangle + 2\omega_{k+1} \langle \mathbf{T}_k, \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) - \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) \rangle \\ &= \|\mathbf{T}_k\|^2 + \text{D1} + \text{D2} + \text{D3}. \end{aligned}$$

First of all, according to (26) and (6), we have

$$\omega_{k+1}^2 \|H(\boldsymbol{\theta}_k, \boldsymbol{\beta}_{k+1})\|^2 \leq \omega_{k+1}^2 C_1 (1 + \|\mathbf{T}_k\|^2), \quad (\text{D1})$$

$$2\omega_{k+1} \langle \mathbf{T}_k, h(\boldsymbol{\theta}_k) \rangle \leq -2\omega_{k+1} \delta \|\mathbf{T}_k\|^2, \quad (\text{D2})$$

Conduct the decomposition of D3 similar to Theorem 24 (p.g. 246) from Benveniste et al. [1990] and Lemma A.5 [Liang, 2010].

$$\begin{aligned} & \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) - \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) \\ &= \underbrace{\mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) - \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_k)}_{\text{D3-1}} + \underbrace{\Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_k) - \Pi_{\boldsymbol{\theta}_{k-1}} \mu_{\boldsymbol{\theta}_{k-1}}(\boldsymbol{\beta}_k)}_{\text{D3-2}} + \underbrace{\Pi_{\boldsymbol{\theta}_{k-1}} \mu_{\boldsymbol{\theta}_{k-1}}(\boldsymbol{\beta}_k) - \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1})}_{\text{D3-3}}. \end{aligned}$$

(i)  $\mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) - \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_k)$  forms a martingale difference sequence such that

$$\mathbb{E} [\mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) - \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_k) | \mathcal{F}_k] = 0. \quad (\text{D3-1})$$

(ii) From Lemma 1, we have that  $\mathbb{E}[\|\mathbf{T}_k\|]$  is bounded.  $\|\Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}\|$  is also bounded according to (25). Therefore, together with Cauchy-Schwarz inequality, there exists a positive constant  $C_2$  such that

$$\mathbb{E} [2\omega_{k+1} \langle \mathbf{T}_k, \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_k) - \Pi_{\boldsymbol{\theta}_{k-1}} \mu_{\boldsymbol{\theta}_{k-1}}(\boldsymbol{\beta}_k) \rangle] \leq \omega_{k+1} C_2. \quad (\text{D3-2})$$

(iii) D3-3 can be further decomposed to D3-3a and D3-3b

$$\begin{aligned} & \langle \mathbf{T}_k, \Pi_{\boldsymbol{\theta}_{k-1}} \mu_{\boldsymbol{\theta}_{k-1}}(\boldsymbol{\beta}_k) - \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) \rangle \\ &= (\langle \mathbf{T}_k, \Pi_{\boldsymbol{\theta}_{k-1}} \mu_{\boldsymbol{\theta}_{k-1}}(\boldsymbol{\beta}_k) \rangle - \langle \mathbf{T}_{k+1}, \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) \rangle) + (\langle \mathbf{T}_{k+1}, \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) \rangle - \langle \mathbf{T}_k, \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) \rangle) \\ &= \underbrace{\langle \mathbf{z}_k - \mathbf{z}_{k+1} \rangle}_{\text{D3-3a}} + \underbrace{\langle \mathbf{T}_{k+1} - \mathbf{T}_k, \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) \rangle}_{\text{D3-3b}}. \end{aligned}$$

where  $\mathbf{z}_k = \langle \mathbf{T}_k, \Pi_{\boldsymbol{\theta}_{k-1}} \mu_{\boldsymbol{\theta}_{k-1}}(\boldsymbol{\beta}_k) \rangle$ . Similar to (ii), there exists a constant  $C_3$  such that

$$\mathbb{E} [2\omega_{k+1} \langle \mathbf{T}_{k+1} - \mathbf{T}_k, \Pi_{\boldsymbol{\theta}_k} \mu_{\boldsymbol{\theta}_k}(\boldsymbol{\beta}_{k+1}) \rangle] \leq C_3 \omega_{k+1}$$

Finally, add all the items D1, D2 and D3 together, for some  $\overline{C}_1 = C_2 + C_3$ , we have

$$\mathbb{E} [\|\mathbf{T}_{k+1}\|^2] \leq (1 - 2\omega_{k+1}\delta + \omega_{k+1}^2 C_1) \mathbb{E} [\|\mathbf{T}_k\|^2] + \omega_{k+1}^2 C_1 + \omega_{k+1} \overline{C}_1 + 2\omega_{k+1} \mathbb{E}[z_k - z_{k+1}].$$

Moreover, from (25), there exists a constant  $C_4$  such that

$$\mathbb{E}[\|z_k\|] \leq C_4. \quad (31)$$

Lemma 4 is an extension of Lemma 26 (page 248) from Benveniste et al. [1990].

**Lemma 4.** *Let  $\{u_k\}_{k \geq k_0}$  as a sequence of real numbers such that for all  $k \geq k_0$ , some suitable constants  $\overline{C}_1$  and  $C_1$*

$$u_{k+1} \geq u_k (1 - 2\omega_{k+1}\delta + \omega_{k+1}^2 C_1) + \omega_{k+1}^2 C_1 + \omega_{k+1} \overline{C}_1, \quad (32)$$

and assume there exists such  $k_0$  that

$$\mathbb{E} [\|\mathbf{T}^{(k_0)}\|^2] \leq u^{(k_0)}. \quad (33)$$

Then for all  $k > k_0$ , we have

$$\mathbb{E} [\|\mathbf{T}_k\|^2] \leq u_k + \sum_{j=k_0+1}^k \Lambda_j^k (\mathbf{z}^{(j-1)} - \mathbf{z}^{(j)}).$$

**Proof of Theorem 1 (Continued).** From Lemma 3, we can choose  $\lambda_0$  and  $k_0$  which satisfy the conditions (32) and (33)

$$\mathbb{E} [\|\mathbf{T}^{(k_0)}\|^2] \leq u^{(k_0)} = \lambda_0 \omega^{(k_0)}.$$

From Lemma 4, it follows that for all  $k > k_0$

$$\mathbb{E} [\|\mathbf{T}_k\|^2] \leq u_k + \mathbb{E} \left[ \sum_{j=k_0+1}^k \Lambda_j^k (\mathbf{z}^{(j-1)} - \mathbf{z}^{(j)}) \right]. \quad (34)$$

From (31) and the increasing property of  $\Lambda_j^k$  in Lemma 2, we have

$$\begin{aligned} & \mathbb{E} \left[ \left\| \sum_{j=k_0+1}^k \Lambda_j^k (\mathbf{z}^{(j-1)} - \mathbf{z}^{(j)}) \right\|^2 \right] \\ &= \mathbb{E} \left[ \left\| \sum_{j=k_0+1}^{k-1} (\Lambda_{j+1}^k - \Lambda_j^k) \mathbf{z}^{(j)} - 2\omega_k \mathbf{z}_k + \Lambda_{k_0+1}^k \mathbf{z}^{(k_0)} \right\|^2 \right] \\ &\leq \sum_{j=k_0+1}^{k-1} (\Lambda_{j+1}^k - \Lambda_j^k) C_4 + \mathbb{E} [2\omega_k \mathbf{z}_k] + \Lambda_{k_0+1}^k C_4 \\ &\leq (\Lambda_k^k - \Lambda_{k_0+1}^k) C_4 + \Lambda_{k_0+1}^k C_4 + \Lambda_{k_0+1}^k C_4 \\ &\leq 3\Lambda_{k_0+1}^k C_4 = 6C_4 \omega_k. \end{aligned} \quad (35)$$

Therefore, given the sequence  $u_k = \lambda_0 \omega_k$  that satisfies conditions (32), (33) and Lemma 4, for any  $k > k_0$ , from (34) and (35), we have

$$\mathbb{E} [\|\mathbf{T}_k\|^2] \leq u_k + 3C_4 \Lambda_k^k = (\lambda_0 + 6C_4) \omega_k = \lambda \omega_k,$$

where  $\lambda = \lambda_0 + 6C_4$ . □

Table 1: Predictive errors in logistic regression based on a test set considering different  $v_0$  and  $\sigma$

MAE / MSE	$v_0=0.01, \sigma=1$	$v_0=0.001, \sigma=1$	$v_0=0.01, \sigma=2$	$v_0=0.001, \sigma=2$
SGLD-SA	<b>0.177 / 0.108</b>	<b>0.188 / 0.114</b>	<b>0.182 / 0.116</b>	<b>0.187 / 0.113</b>
SGLD-EM	0.207 / 0.131	0.361 / 0.346	0.204 / 0.132	0.376 / 0.360
SGLD	0.295 / 0.272	0.335 / 0.301	0.350 / 0.338	0.337 / 0.319

## 2.2 Weak Convergence of Samples

In statistical models with latent variables, the gradient is often biased due to the use of stochastic approximation. Langevin Monte Carlo with inaccurate gradients has been studied by Chen et al. [2015], Dalalyan and Karagulyan [2018], which are helpful to prove the weak convergence of samples. Following theorem 2 in Chen et al. [2015], we have

**Corollary 1.** *Under assumptions in Appendix B.1 and the assumption 1 (smoothness and boundness on the solution functional) in Chen et al. [2015], the distribution of  $\beta_k$  converges weakly to the target posterior as  $\epsilon \rightarrow 0$  and  $k \rightarrow \infty$ .*

*Proof.* Since  $\theta_k$  converges to  $\theta^*$  in SGLD-SA under assumptions in Appendix B.1 and the gradient is M-smooth (9), we decompose the stochastic gradient  $\nabla_{\beta} \tilde{L}(\beta_k, \theta_k)$  as  $\nabla_{\beta} L(\beta_k, \theta^*) + \xi_k + \mathcal{O}(k^{-\alpha})$ , where  $\nabla_{\beta} L(\beta_k, \theta^*)$  is the exact gradient,  $\xi_k$  is a zero-mean random vector,  $\mathcal{O}(k^{-\alpha})$  is the bias term coming from the stochastic approximation and  $\alpha \in (0, 1]$  is used to guarantee the consistency in theorem 1. Therefore, Eq.(3) can be written as

$$\beta_{k+1} = \beta_k + \epsilon_k (\nabla_{\beta} L(\beta_k, \theta^*) + \xi_k + \mathcal{O}(k^{-\alpha})) + \sqrt{2\epsilon_k} \eta_k, \text{ where } \eta_k \sim \mathcal{N}(0, I). \quad (36)$$

Following a similar proof in Chen et al. [2015], it suffices to show that  $\sum_{k=1}^K k^{-\alpha} / K \rightarrow 0$  as  $K \rightarrow \infty$ , which is obvious. Therefore, the distribution of  $\beta_k$  converges weakly to the target distribution as  $\epsilon \rightarrow 0$  and  $k \rightarrow \infty$ .  $\square$

## 3 Simulation of Large-p-Small-n Logistic Regression

Now we conduct the experiments on binary logistic regression. The setup is similar as before, except  $n$  is set to 500,  $\Sigma_{i,j} = 0.3^{|i-j|}$  and  $\eta \sim \mathcal{N}(0, I/2)$ . We set the learning rate for all the three algorithms to  $0.001 \times k^{-\frac{1}{3}}$  and step size  $\omega_k$  to  $10 \times (k + 1000)^{-0.7}$ . The binary response values are simulated from **Bernoulli**( $p$ ) where  $p = 1/(1 + e^{-X\beta - \eta})$ . As shown in Fig.1: SGLD fails in selecting the right variables and overfits the data; both SGLD-EM and SGLD-SA choose the right variables. However, SGLD-EM converges to a poor local optimum by mistakenly using  $L_1$  norm to regularize all the variables, leading to a large shrinkage effect on  $\beta_{1:3}$ . By contrast, SGLD-SA successfully updates the latent variables and regularize  $\beta_{1:3}$  with  $L_2$  norm, yielding a better parameter estimation for  $\beta_{1:3}$  and a stronger regularization for  $\beta_{4-1000}$ . Table.1 illustrates that SGLD-SA consistently outperforms the other methods and is robust to different initializations. We observe that SGLD-EM sometimes performs as worse as SGLD when  $v_0 = 0.001$ , which indicates that the EM-based variable selection is not robust in the stochastic optimization of the latent variables.

## 4 Regression on UCI datasets

We further evaluate our model on five **UCI** regression datasets and show the results in Table 2. Following Hernandez-Lobato and Adams [2015], we randomly sample 90% of each dataset for training and leave the rest for testing. We run 20 experiments for each setup with fixed random seeds and report the averaged error rate. Feature normalization is applied in the experiments. The model is a simple MLP with one hidden layer of 50 units. We set the batch size to 50, the training epoch to 200, the learning rate to  $1e-5$  and the default  $L_2$  to  $1e-4$ . For SGHMC-EM and SGHMC-SA, we apply the SSGL prior on the BNN weights (excluding biases) and fix  $a, \nu, \lambda = 1, b, v_1, \sigma = 10$  and  $\delta = 0.5$ . We fine-tune the initial temperature  $\tau$  and  $v_0$ . As shown in Table 2, SGHMC-SA outperforms all the baselines. Nevertheless, without smooth adaptive update, SGHMC-EM often performs worse than SGHMC. While with simulated annealing where  $\tau^{(k)} = \tau \times 1.003^k$ , we observe further improved performance in most of the cases.



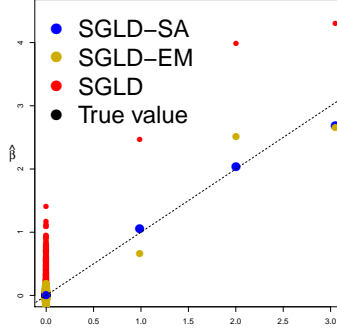


Figure 1: Logistic regression simulation when  $v_0 = 0.1$  and  $\sigma = 1$

Dataset Hyperparameters	Boston 1/0.1	Yacht 1/0.1	Energy 0.1/0.1	Wine 0.5/0.01	Concrete 0.5/0.07
SGHMC	2.783±0.109	0.886±0.046	1.983±0.092	0.731±0.015	6.319±0.179
A-SGHMC	2.848±0.126	0.808±0.048	1.419±0.067	0.671±0.019	5.978±0.166
SGHMC-EM	2.813±0.140	0.823±0.053	2.077±0.108	0.729±0.018	6.275±0.169
A-SGHMC-EM	2.767±0.154	0.815±0.052	1.435±0.069	0.627±0.008	5.762±0.156
SGHMC-SA	<b>2.779±0.133</b>	<b>0.789±0.050</b>	<b>1.948±0.081</b>	<b>0.654±0.010</b>	<b>6.029±0.131</b>
A-SGHMC-SA	<b>2.692±0.120</b>	<b>0.782±0.052</b>	<b>1.388±0.052</b>	<b>0.620±0.008</b>	<b>5.687±0.142</b>

Table 2: Average performance and standard deviation of Root Mean Square Error, where  $\tau$  denotes the initial inverse temperature and  $v_0$  is a hyperparameter in the SSGL prior (Hyperparameters  $\tau/v_0$ ).

## 5 Experimental Setup

### 5.1 Network Architecture

The first DNN we use is a standard 2-Conv-2-FC CNN: it has two convolutional layers with a  $2 \times 2$  max pooling after each layer and two fully-connected layers. The filter size in the convolutional layers is  $5 \times 5$  and the feature maps are set to be 32 and 64, respectively [Jarrett et al., 2009]. The fully-connected layers (FC) have 200 hidden nodes and 10 outputs. We use the rectified linear unit (ReLU) as activation function between layers and employ a cross-entropy loss.

The second DNN is a 2-Conv-BN-3-FC CNN: it has two convolutional layers with a  $2 \times 2$  max pooling after each layer and three fully-connected layers with batch normalization applied to the first FC layer. The filter size in the convolutional layers is  $4 \times 4$  and the feature maps are both set to 64. We use  $256 \times 64 \times 10$  fully-connected layers.

### 5.2 Data Augmentation

The MNIST dataset is augmented by (1) randomCrop: randomly crop each image with size 28 and padding 4, (2) random rotation: randomly rotate each image by a degree in  $[-15^\circ, +15^\circ]$ , (3) normalization: normalize each image with empirical mean 0.1307 and standard deviation 0.3081.

The FMNIST dataset is augmented by (1) randomCrop: same as MNIST, (2) randomHorizontalFlip: randomly flip each image horizontally, (3) normalization: same as MNIST, (4) random erasing [Zhong et al., 2017].

The CIFAR10 dataset is augmented by (1) randomCrop: randomly crop each image with size 32 and padding 4, (2) randomHorizontalFlip: randomly flip each image horizontally, (3) normalization: normalize each image with empirical mean (0.4914, 0.4822, 0.4465) and standard deviation (0.2023, 0.1994, 0.2010), (4) random erasing.

## References

- Albert Benveniste, Michael Métivier, and Pierre Priouret. *Adaptive Algorithms and Stochastic Approximations*. Berlin: Springer, 1990.
- Changyou Chen, Nan Ding, and Lawrence Carin. On the Convergence of Stochastic Gradient MCMC Algorithms with High-order Integrators. In *Proc. of the Conference on Advances in Neural Information Processing Systems (NIPS)*, pages 2278–2286, 2015.
- Arnak S. Dalalyan and Avetik G. Karagulyan. User-friendly guarantees for the Langevin Monte Carlo with inaccurate gradient. *ArXiv e-prints*, September 2018.
- Jose Miguel Hernandez-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *Proc. of the International Conference on Machine Learning (ICML)*, volume 37, pages 1861–1869, 2015.
- K. Jarrett, K. Kavukcuoglu, M. Ranzato, and Y. LeCun. What is the best multi-stage architecture for object recognition? In *Proc. of the International Conference on Computer Vision (ICCV)*, pages 2146–2153, September 2009.
- Faming Liang. Trajectory averaging for stochastic approximation MCMC algorithms. *The Annals of Statistics*, 38:2823–2856, 2010.
- Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random Erasing Data Augmentation. *ArXiv e-prints*, 2017.