Machine Learning and Neural Networks (MATH3431)

Epiphany term, 2025

Lecture notes 7: Bayesian Learning via Stochastic gradient and Stochastic gradient Langevin dynamics

Lecturer & author: Georgios P. Karagiannis georgios.karagiannis@durham.ac.uk

Aim. To introduce the Bayesian Learning, and Stochastic gradient Langevin dynamics (description, heuristics, and implementation). Stochastic gradient descent will be implemented in the Bayesian learning too.

Reading list & references:

- (1) Welling, M., & Teh, Y. W. (2011). Bayesian learning via stochastic gradient Langevin dynamics. In Proceedings of the 28th international conference on machine learning (ICML-11) (pp. 681-688).
- (2) Bottou, L. (2012). Stochastic gradient descent tricks. In Neural networks: Tricks of the trade (pp. 421-436). Springer, Berlin, Heidelberg.
- (3) Sato, I., & Nakagawa, H. (2014). Approximation analysis of stochastic gradient Langevin dynamics by using Fokker-Planck equation and Ito process. In International Conference on Machine Learning (pp. 982-990). PMLR.
- (4) Teh, Y. W., Thiery, A. H., & Vollmer, S. J. (2016). Consistency and fluctuations for stochastic gradient Langevin dynamics. Journal of Machine Learning Research, 17.
- (5) Vollmer, S. J., Zygalakis, K. C., & Teh, Y. W. (2016). Exploration of the (non-) asymptotic bias and variance of stochastic gradient Langevin dynamics. The Journal of Machine Learning Research, 17(1), 5504-5548. ->further reading for theory
- (6) Nemeth, C., & Fearnhead, P. (2021). Stochastic gradient Markov chain Monte Carlo. Journal of the American Statistical Association, 116(533), 433-450.

1. Bayesian learning and motivations

Note 1. Bayesian methods are appealing in their ability to capture uncertainty in learned parameters and avoid overfitting. Arguably with large datasets there will be little overfitting. Alternatively, as we have access to larger datasets and more computational resources, we become interested in building more complex models (eg from a logistic regression to a deep neural network), so that there will always be a need to quantify the amount of parameter uncertainty.

Note 2. Consider a Bayesian statistical model with sampling distribution (statistical model) f(z|w) labeled by an unknown parameter $w \in \mathcal{W} \subseteq \mathbb{R}^d$ that follows a prior distribution f(w). Assume a dataset $\mathcal{S}_n = \{z_i; i = 1, ..., n\}$ of size n containing independently drawn examples. Let $L_n(w) := f(z_{1:n}|w)$ denote the likelihood of the observables $\{z_i \in \mathcal{Z}\}_{i=1}^n$ give parameter w. The Bayesian

model is denoted as

(1.1)
$$\begin{cases} z_i | w \stackrel{\text{ind}}{\sim} f(z_i | w), \ i = 1, ..., n \\ w \sim f(w) \end{cases}$$

Note 3. With regards to the Bayesian model in Note 2, we denote the likelihood of the observables $\{z_i \in \mathcal{Z}\}_{i=1}^n$ given the parameter w as

$$L_n(w) := f(z_{1:n}|w) = \prod_{i=1}^n f(z_i|w)$$

Note 4. Bayesian learning (inference) relies on the posterior distribution density

(1.2)
$$f(w|z_{1:n}) = \frac{L_n(w) f(w)}{\int L_n(w) f(w) dw}$$

which quantifies the researcher's belief (or uncertainty) about the unknown parameter w learned given examples $\{z_i \in \mathcal{Z}\}_{i=1}^n$. It is often intractable; hence there is often a need to numerically compute it. (Section 3)

Note 5. Point estimation of a function h of w is often performed via computation of the posterior expectation w given the examples in S_n

(1.3)
$$E_f(h(w)|z_{1:n}) = \int h(w) f(w|z_{1:n}) dw$$

It is often intractable; hence there is often a need to numerically compute it. (Section 3)

Note 6. Point estimation of w can also be performed via maximum a-posteriori (MAP) estimator w^* of w that is the maximizer w^* of (1.2) i.e.

(1.4)
$$w^* = \arg\max_{w} (f(w|z_{1:n}))$$

(1.5)
$$= \underset{w}{\operatorname{arg\,max}} \left(\underbrace{\log\left(L_{n}\left(w\right)\right)}_{\text{(I)}} + \underbrace{\log\left(f\left(w\right)\right)}_{\text{(II)}} \right)$$

Note that (I) may be interpreted as an empirical risk function, and (II) can be interpreted as a shrinkage term in terms of shrinkage methods (like LASSO, Ridge). It is often intractable; hence there is often a need to numerically compute it. (Section 2)

Note 7. To describe the learning algorithms Gradient Descent (GD), Stochastic Gradient Descent (SGD), and Stochastic Gradient Langevin Dynamics (SGLD), we consider that the examples (data) in (1.1) are independent realizations from the sampling distribution i.e. $z_i|w \sim f(\cdot|w)$ for i = 1, ..., n and that w is continuous (not discrete).

Note 8. In what follows, we first present the implementation of GD and SGD addressing MAP learning, and then we introduce the implementation of SGLD addressing posterior density and expectation learning.

2. Maximum A Posteriori (MAP) learning via GD and SGD

Problem 9. Given the Bayesian model (1.1), and rearranging (1.4), MAP estimate w^* of w can be computed as

(2.1)
$$w^* = \underset{w}{\operatorname{arg max}} \left(\log \left(L_n(w) \right) + f(w) \right) = \underset{w}{\operatorname{arg max}} \left(\sum_{i=1}^n \log \left(f(z_i | w) \right) + \log \left(f(w) \right) \right)$$

Note 10. GD is particularly suitable to solve (2.1) when w has high dimensionality.

Algorithm 11. Gradient descent implementation in Problem 9 with learning rate $\eta_t \geq 0$.

For $t = 1, 2, 3, \dots$ iterate:

(1) Compute

$$(2.2) w^{(t+1)} = w^{(t)} + \eta_t \left(\sum_{i=1}^n \nabla_w \log \left(f\left(z_j | w^{(t)}\right) \right) + \nabla_w \log \left(f\left(w^{(t)}\right) \right) \right)$$

Note 12. SGD is particularly suitable to solve (2.1) when w has high dimensionality, and in bigdata problems since the repetitive computation of the sum in (2.2) is prohibitively expensive. Yet consider the benefits of SGD against GD as discussed in (Notes 27 and 28 of Lect. Notes 5).

Algorithm 13. Stochastic Gradient Descent implementation in with learning rate $\eta_t \geq 0$ and batch size m

For t = 1, 2, 3, ... iterate:

- (1) generate a random set $\mathcal{J}^{(t)} \subseteq \{1,...,n\}^m$ of m indices from 1 to n with or without replacement
- (2) compute

(2.3)
$$w^{(t+1)} = w^{(t)} + \eta_t \left(\frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \log \left(f\left(z_j | w^{(t)}\right) \right) + \nabla_w \log \left(f\left(w^{(t)}\right) \right) \right)$$

Note 14. Recursion (2.3) is justified in terms of SGD theory as

$$(2.4) \quad \mathrm{E}_{\mathcal{J}^{(t)} \sim \mathrm{simple-random-sampling}} \left(\frac{n}{m} \sum_{j \in \mathcal{J}^{(t)}} \nabla_w \log \left(f\left(z_j | w^{(t)}\right) \right) \right) = \sum_{i=1}^n \nabla_w \log \left(f\left(z_i | w^{(t)}\right) \right)$$

see Exercise 15 from the Exercise sheet.

Note 15. The implementation of the SGD variants (Algorithms 23, 25, 10, 15 in Handout 3)) is straightforward based on Algorithm 13 and (2.4).

3. Fully Bayesian learning via SGLD

Problem 16. Fully Bayesian learning, computationally, is the problem of recovering the posterior distribution $f(w|z_{1:n})$ of w given $z_{1:n}$ that admits density (1.2). For a given Bayesian model (1.1), Page 3 Created on 2025/02/10 at 11:12:27 by Georgios Karagiannis

the Bayesian estimator of h := h(w) is often computed as the posterior expectation of h(w) given the data S_n

(1.3)
$$E_f(h(w)|z_{1:n}) = \int h(w) f(w|z_{1:n}) dw$$

Note 17. Monte Carlo integration aims at approximating (1.3) when intractable, by using Central Limit Theorem or Law of Large Numbers arguments as $h_T \approx E_f(h(w)|z_{1:n})$ where

(3.1)
$$\hat{h}_T = \frac{1}{T} \sum_{t=1}^{T} h\left(w^{(t)}\right)$$

where $\{w^{(t)}\}\$ are T simulations drawn (approximately) from the posterior distribution (1.2). This theory is subject to conditions we skip.

Note 18. Stochastic gradient Langevin dynamics (SGLD) algorithm generates as output a random chain $\{w^{(t)}\}$ approximately distributed according to a distribution admitting density such as

(3.2)
$$f_{\tau}\left(w|z_{1:n}\right) \propto \exp\left(\frac{1}{\tau}\left(\log L_{n}\left(w\right) + \log f\left(w\right)\right)\right)$$

(3.3)
$$\propto \exp\left(\frac{1}{\tau} \left(\sum_{i=1}^{n} \log\left(f\left(z_{i}|w\right)\right) + \log\left(f\left(w\right)\right)\right)\right)$$

for any $\tau > 0$ under regularity conditions. That allows to recover the whole posterior distribution (1.2) (hence account for uncertainty in h = h(w)) and approximate posterior expectations based on the Monte Carlo integration (Note 17).

Algorithm 19. Stochastic Gradient Langevin Dynamics (SGLD) with learning rate $\eta_t > 0$, batch size m, and temperature $\tau > 0$ is

- Ref [1]
- (1) Generate a random set $\mathcal{J}^{(t)} \subseteq \{1,...,n\}^m$ of m indices from 1 to n with or without replacement
- (2) Compute

$$(3.4) w^{(t+1)} = w^{(t)} + \eta_t \left(\frac{n}{m} \sum_{i \in J^{(t)}} \nabla \log f\left(z_i | w^{(t)}\right) + \nabla \log f\left(w^{(t)}\right) \right) + \sqrt{2\eta_t \tau} \epsilon_t$$

where $\epsilon_t \stackrel{iid}{\sim} N(0, I)$.

(3) Terminate if a termination criterion is satisfied; E.g., $t \geq T_{max}$ for a prespecified $T_{max} >$

How / why it works...

Note 20. (Mathematically speaking) The stochastic chain in (3.4) can be viewed as a discretization Ref [3] of the continuous-time Langevin diffusion described by the stochastic differential equation

(3.5)
$$dW(t) = -\nabla_w \left[-\log f(W(t)|z_{1:n}) \right] dt + \sqrt{2\tau} dB(t), \ t \ge 0$$

Created on 2025/02/10 at 11:12:27

by Georgios Karagiannis

where $\{B(t)\}\$ is a standard Brownian motion in \mathbb{R}^d (i.e.). Under suitable assumptions on f, it can be shown that a Gibbs distribution with PDF such as

(3.6)
$$f^*(w|z_{1:n}) \propto \exp\left(-\frac{1}{\tau} \left[-\log f(w|z_{1:n})\right]\right)$$

is the unique invariant distribution of (3.5), and that the distributions of W(t) converge rapidly to

Note 21. (Heuristically speaking) SGLD relies on injecting the 'right' amount of noise (3.4) to a standard stochastic gradient optimization recursion (2.2), such that, as the stepsize η_t properly reduces, the produced chain $\{w^{(t+1)}\}$ converges to samples that could have been drawn from distribution with density (3.2). In the initial phase of running, the stochastic gradient noise will dominate the injected noise ϵ_t and the algorithm will imitate an efficient SGD Algorithm 11 -but this is until η_t or $\nabla \log (L_n(w))$ become small enough. In the later phase of running, the injected noise ϵ_t will dominate the stochastic gradient noise, so the SGLD will imitate a Langevin dynamics for the target distribution (1.2). The aim is for to the algorithm to transition smoothly between the two phases. Whether the algorithm is in the stochastic optimization phase or Langevin dynamics phase depends on the variance of the injected noise versus that of the stochastic gradient.

Note 22. One can argue that, the output of SGLD is also an "almost" maximizer of the (minus) empirical risk with shrinkage for large enough t. A draw from the Gibbs distribution (3.6) is approximately a maximizer of (2.1). Also one can show that the SGLD recursion tracks the Langevin diffusion (3.5) in a suitable sense. Hence, both imply that the distributions of W(t) will be close to the Gibbs distribution (3.6) for all sufficiently large t.

About the learning rate alg. parameter...

Note 23. To guarantee the algorithm to work it is important for the step sizes η_t to decrease to zero, so that the mixing rate of the algorithm will slow down with increasing number of iterations t. Then, we can keep the step size η_t constant once it has decreased below a critical level.

Condition 24. Regarding the learning rate (or gain) $\{\eta_t\}$ should satisfy the following conditions in order for SGLD Algorithm 19 to generate an output $\{w^{(t)}\}$ approximately distributed according to (3.2)

- (1) $\eta_t \ge 0$,
- $(2) \sum_{t=1}^{\infty} \eta_t = \infty$ $(3) \sum_{t=1}^{\infty} \eta_t^2 < \infty$

Note 25. The popular learning rates $\{\eta_t\}$ in Note 12 in Lect. Notes 4 can be used. . Once parametrized, η_t can be tuned based on pilot runs using a reasonably small number of data.

 $^{^{1}}$ A a continuous-time stochastic process: (1) B(0) = 0; (2) B(t) is almost surely continuous; (3) B(t) has independent increaments; (4) $B(t) - B(s) \sim N(0, t - s)$ for $0 \le s \le t$.

About the temperature alg. parameter...

Note 26. The temperature parameter $\tau > 0$ is user specified and aims at controlling (eg; inflating) the variance of the produced chain for instance with practical purpose to escape from local modes (otherwise energy barriers) in non-convex problems.

About the output of the alg...

Note 27. The first few iterations of SGLD (Algorithm 19) involve values generated at the beginning of the running algorithm while the chain have not yet converged to (or reached) an area of substantial posterior mass. Hence they are discarded from the output of the SGLD. These values are called burn-in.

Note 28. The output of SGLD (Algorithm 19) $\{w^{(t)}\}$ includes the generated values of w produced during the last few iterations of the running algorithm and after discarding the burn-in generated values.

Note 29. SGLD for $\tau = 1$ approximately simulates from the posterior (1.2).

About facilitating inference...

Note 30. Expectation (1.3), can be estimated as an arithmetic average

(3.7)
$$\widehat{h_T(w)} = \frac{1}{T} \sum_{t=1}^{T} h\left(w^{(t)}\right)$$

as $\widehat{h_T(w)} \to \mathrm{E}_f(h(w)|z_{1:n})$ based on LLN arguments.

Note 31. Another more efficient estimator for the expectation (1.3) is the weighted arithmetic average

$$\widehat{\hat{h}_T(w)} = \sum_{t=T_0+1}^T \frac{\eta_t}{\sum_{t=T_0+1}^T \eta_{t'}} h\left(w^{(t)}\right)$$

Note 32. Estimator (3.8) has a smaller standard error than estimator (3.7). As the step size η_t decreases, the mixing rate of the chain $\{w^{(t)}\}$ decreases as well and the simple sample average (3.7) overemphasizes in the tail end of the sequence where there is higher correlation among the samples resulting in higher variance in the estimator.

About the exploding gradient phenomenon...

Note 33. 'Exploding gradients' is the practical phenomenon in which large updates to weights during training can cause a numerical overflow or underflow due to the machine error of the computer.

Note 34. Practical solutions to 'Exploding gradient' involve, at each iteration t, checking the magnitude of the gradient v_t , (e.g., Euclidean norm $||v_t||$), and instantly changing it (e.g., gradient scaling or clipping) if it is gonna result an overflow.

Note 35. Gradient scaling involves normalizing the gradient vector such that vector norm (magnitude) equals a user-specified value. More formally, given a gradient v_t , gradient clipping can be used in a standard recursion

$$w^{(t+1)} = w^{(t)} + \eta_t v_t$$

as

$$w^{(t+1)} = w^{(t)} + \eta_t \text{clip}(v_t, c)$$

where

$$\operatorname{clip}(v,c) = v \min\left(1, \frac{c}{\|v\|}\right)$$

and c > 0 is a clipping threshold impaling that clipping will take place at iteration t if $||v_t|| > c$.

Note 36. Gradient clipping involves forcing the gradient values (element-wise) to a specific minimum or maximum value if the gradient exceeded an expected range.

Note 37. Unreasonably frequent scaling or clipping may introduce significant bias and hence it should not be applied carelessly.

We continue the Example 26 in 4: Gradient descent.

Specifying the Bayesian model. Consider the Bayesian Normal linear regression model

$$\begin{cases} y_i | \beta, \sigma^2 \sim \mathrm{N}\left(x_i^{\top} \beta, \sigma^2\right) & \text{sampling distribution } f\left(y_i | \beta, \sigma^2\right) \\ \beta | \sigma^2 \sim \mathrm{N}\left(\mu, \sigma^2 V\right) & \text{prior } f\left(\beta | \sigma^2\right) \\ \sigma^2 \sim \mathrm{IG}\left(\phi, \psi\right) & \text{prior } f\left(\sigma^2\right) \end{cases}$$

and $f(\beta, \sigma^2) = f(\beta | \sigma^2) f(\sigma^2)$, $\beta \in \mathbb{R}^d$, and $\sigma^2 \in \mathbb{R}_+$. Note given densities

$$\begin{split} \mathbf{N}\left(x|\mu,\Sigma\right) &= \left(\frac{1}{2\pi}\right)^{d} \frac{1}{|\Sigma|} \exp\left(-\frac{1}{2}\left(x-\mu\right)^{\top} \Sigma^{-1}\left(x-\mu\right)\right) \\ \mathbf{IG}\left(x|a,b\right) &= \frac{b^{a}}{\Gamma\left(a\right)} x^{-a-1} \exp\left(-\frac{b}{x}\right) \mathbf{1}\left(x \geq 0\right) \end{split}$$

Performing computationally convenient transformations. Because SGD (Algorithm 13) and SGLD (Algorithm 19) can handle cases where $w \in \mathbb{R}^d$ in a straightforward manner than what they do when $w = (\beta, \sigma^2) \in \mathbb{R}^d \times \mathbb{R}_+$ which requires an additional projection step; we consider a transformation $w = (\beta, \gamma) \in \mathbb{R}^d \times \mathbb{R}$ with $\gamma = \log(\sigma^2)$. Hence, the Bayesian model becomes

$$\begin{cases} y_{i}|\beta, \sigma^{2} \sim \mathcal{N}\left(x_{i}^{\top}\beta, \exp\left(\gamma\right)\right) & \text{sampling distribution } f\left(y_{i}|\beta, \gamma\right) \\ \beta|\sigma^{2} \sim \mathcal{N}\left(\mu = 0, \exp\left(\gamma\right)V\right) & \text{prior } f\left(\beta|\gamma\right) \\ \gamma \sim f_{\gamma}\left(\gamma\right) & \text{prior } f\left(\gamma\right) \end{cases}$$

²Code is available from https://github.com/georgios-stats/Machine_Learning_and_Neural_Networks_III_ Epiphany_2025/tree/main/Lecture_notes/code/06.Stochastic_gradient_Langevine_dynamics

where $f_{\gamma}(\gamma)$ is computed according to the method of bijective transformation of random variables; i.e.

$$f_{\gamma}(\gamma) = \operatorname{IG}(\exp(\gamma) | \phi, \psi) \left| \frac{\mathrm{d}}{\mathrm{d}\gamma} \exp(\gamma) \right| = \operatorname{IG}(\exp(\gamma) | \phi, \psi) \exp(\gamma)$$

Computing the require quantities for SGLD. Then we can compute the required gradients in order to run the SGD, and SGLD with respect to $w = (\beta, \gamma)$ with $\gamma = \log(\sigma^2)$.

$$\log \left(f\left(z_i = \left(x_i, y_i \right) | w \right) \right) = -\frac{1}{2} \log \left(2\pi \right) - \frac{1}{2} \gamma - \frac{1}{2} \left(y_j - x_i^\top \beta \right)^2 \exp \left(-\gamma \right)$$

$$\log \left(f\left(w = \left(\beta, \gamma \right) \right) \right) = \log \left(f\left(\beta | \gamma \right) \right) + \log \left(f\left(\gamma \right) \right)$$

$$\log \left(f\left(\beta | \gamma \right) \right) = -\frac{d}{2} \log \left(2\pi \right) - \frac{d}{2} \gamma - \frac{1}{2} \left| V \right| - \frac{1}{2} \exp \left(-\gamma \right) \left(\beta - \mu \right)^\top V^{-1} \left(\beta - \mu \right)$$

$$\log \left(f\left(\gamma \right) \right) = \psi \log \left(\phi \right) - \log \left(\Gamma \left(\phi \right) \right) - \left(\phi + 1 \right) \gamma - \psi \exp \left(-\gamma \right) + \gamma$$

Hence for the log sampling PDF we have

$$\nabla_{w} \log (f(z_{i}|w)) = \left(\frac{\mathrm{d}}{\mathrm{d}\beta} \log (f(z_{i}|w)); \frac{\mathrm{d}}{\mathrm{d}\gamma} \log (f(z_{i}|w))\right)$$
$$\frac{\mathrm{d}}{\mathrm{d}\beta} \log (f(z_{i}|w)) = \left(y_{i} - x_{i}^{\top}\beta\right) x_{i} \exp(-\gamma)$$
$$\frac{\mathrm{d}}{\mathrm{d}\gamma} \log (f(z_{i}|w)) = -\frac{1}{2} + \frac{1}{2} \left(y_{j} - x_{i}^{\top}\beta\right)^{2} \exp(-\gamma)$$

 \dots so

(4.1)
$$\nabla_w \log \left(f\left(z_i | w\right) \right) = \begin{pmatrix} \left(y_i - x_i^\top \beta \right) x_i \exp\left(-\gamma\right) \\ -\frac{1}{2} + \frac{1}{2} \left(y_j - x_i^\top \beta \right)^2 \exp\left(-\gamma\right) \end{pmatrix}$$

Hence for the log a priori PDF we have

$$\nabla_{w} \log (f(w)) = \left(\frac{\mathrm{d}}{\mathrm{d}\beta} \log (f(w)); \frac{\mathrm{d}}{\mathrm{d}\gamma} \log (f(w))\right)$$

$$\frac{\mathrm{d}}{\mathrm{d}\beta} \log (f(w)) = -\exp(-\gamma) V^{-1} (\beta - \mu)$$

$$\frac{\mathrm{d}}{\mathrm{d}\gamma} \log (f(w)) = -\frac{d}{2} + \frac{1}{2} \exp(-\gamma) (\beta - \mu)^{\top} V^{-1} (\beta - \mu) - (\phi + 1) + \psi \exp(-\gamma) + 1$$

...so

(4.2)
$$\nabla_{w} \log (f(w)) = \begin{pmatrix} -\exp(-\gamma) V^{-1} (\beta - \mu) \\ \frac{d}{2} + \frac{1}{2} \exp(-\gamma) (\beta - \mu)^{\top} V^{-1} (\beta - \mu) - (\phi + 1) + \psi \exp(-\gamma) + 1 \end{pmatrix}$$

Implementation of SGLD. We consider $\mu = 0$, $\phi = 1$, $\psi = 1$, and $V = 100I_d$, for our simulations below.

To implement SGD (Algorithm 13) and SGLD (Algorithm 19), we just need to plug in the computed gradients (4.1) and (4.2) for $w = (\beta, \gamma) \in \mathbb{R}^d \times \mathbb{R}$ with $\gamma = \log (\sigma^2)$. After running SGD and SGLD with the computed gradients with respect to $w = (\beta, \gamma) \in \mathbb{R}^d \times \mathbb{R}$ with $\gamma = \log (\sigma^2)$, and obtaining chains $\{w^{(t)} = (\beta^{(t)}, \gamma^{(t)})\}_{t=1}^T$, we can just perform transformation $\{(\sigma^{(t)})^2 = \exp(\gamma^{(t)})\}_{t=1}^T$ if we are interested in learning (β, σ^2) .

In Figures 4.1, we ran the SGD for different batch sizes m and compared it against the exact MLE. We observe that SGD trace converges to the exact MLE. The oscillations are due to the stochastic gradient (ie, noise in the gradient).

In Figures 4.2, we ran the SGLD for different batch sizes m but the same temperature $\tau=1$ and compared it against the exact posterior densities. We observe that the histograms of $\{w^{(t)}\}$ produced from SGLD are closer to the curves representing the exact posteriors when the batch size m is bigger. As we said this is not a panacea; if the landscape of the exact psterior density was multimodal (aka not convex but with had several maxima), then the SGLD using smaller batch sizes could have performed better, in the sense that the inflated noise from the stochastic gradient could accidentally make the generated chain to pass the low mass barrier and visit a different mode, unlike the one with larger batch-size and hence smaller variation.

In Figures 4.3, we ran the SGLD for different temperatures τ but the same batch sizes m=100 and compared it against the exact posterior densities. We observe that increasing the temperature τ may increase the variation of the produced chain. We can use a large τ at the beginning of the run of the algorithm to perform an exploration of the space (this is particularly useful for non-convex/multimodal densities as it allows visiting different modes), and later on we can use a smaller temperature such as $\tau=1$.



FIGURE 4.1. Bayesian learning via SGD (SGD vs (exact)MLE) Study on batch size m.

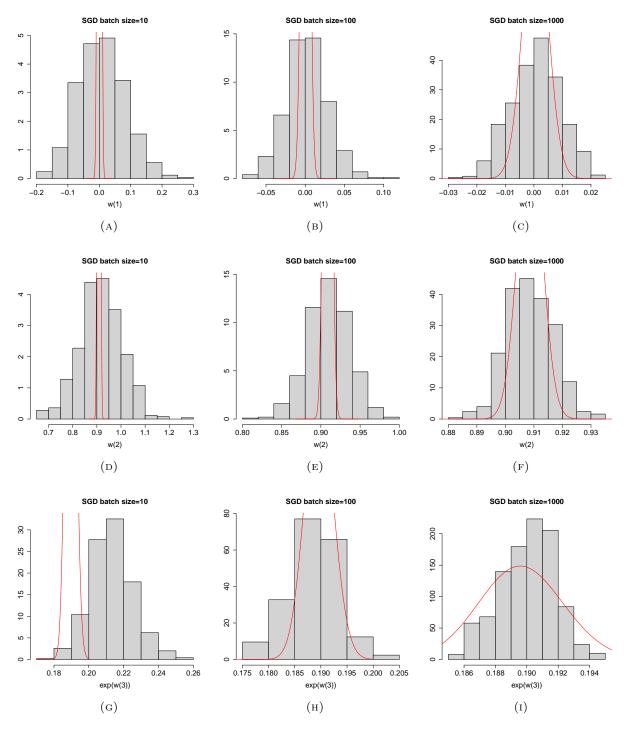


FIGURE 4.2. Bayesian learning: SGLD vs exact posterior (in red). Temperature $\tau=1$. Study on batch size m

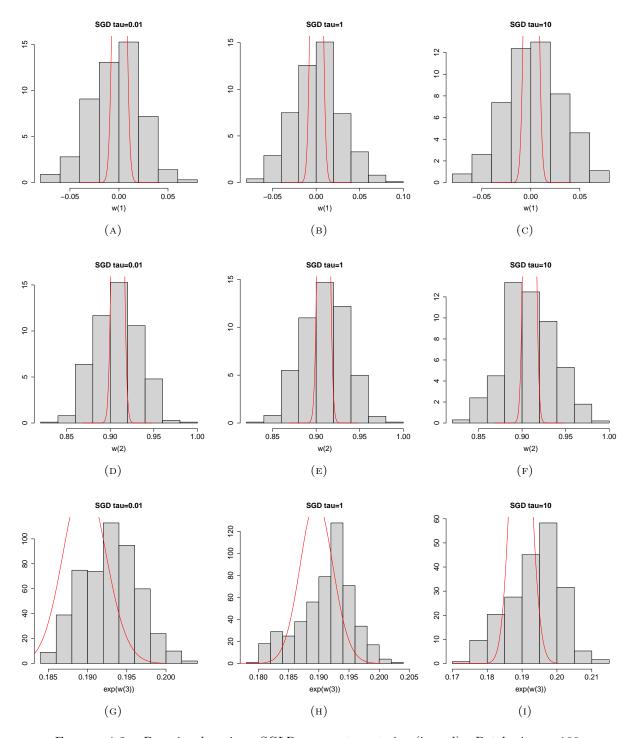


FIGURE 4.3. Bayesian learning: SGLD vs exact posterior (in red). Batch size = 100. Study on τ