## Handout 4: 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:

- 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).
- Bottou, L. (2012). Stochastic gradient descent tricks. In Neural networks: Tricks of the trade (pp. 421-436). Springer, Berlin, Heidelberg.
- 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

## 1. Bayesian learning and motivations

Remark 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  $w \in \Theta \subseteq \mathbb{R}^d$  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) 
$$\begin{cases} z_{1:n}|w & \sim f(z_{1:n}|w) \\ w & \sim f(w) \end{cases}$$

Remark 3. 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)

Remark 4. 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)

Remark 5. 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.

$$(1.4) w^* = \arg\max_{\forall w \in \Theta} \left( f\left(w|z_{1:n}\right) \right)$$

(1.5) 
$$= \arg\min_{\forall w \in \Theta} \left( \underbrace{-\log(L_n(w))}_{(I)} - \underbrace{\log(f(w))}_{(II)} \right)$$

Note that (I) may be interpreted as an empirical rist 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 6. 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* 7. 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 8.** Given the Bayesian model (1.1), and rearranging (1.4), MAP estimate  $w^*$  of w can be computed as

$$(2.1) w^* = \arg\min_{\forall w \in \Theta} \left(-\log\left(L_n\left(w\right)\right) - f\left(w\right)\right) = \arg\min_{\forall w \in \Theta} \left(-\sum_{i=1}^n \log\left(f\left(w|z_i\right)\right) - \log\left(f\left(w\right)\right)\right)$$

Remark 9. GD is particularly suitable to solve (2.1) when w has high dimensionality.

**Algorithm 10.** Gradient descent (Algorithm 1 Handout 2) with learning rate  $\eta_t \geq 0$  can be used to solve (2.1) by using the update rule as

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

(1) Compute

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

Remark 11. The implementation of other GD variants (eg (3.2) in Handout 2) is straightforward, based on Algorithm 10.

Remark 12. SGD is particularly suitable to solve (2.1) when w has high dimensionality, and in big-data 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 Remarks 32 and 33 of Handout 3.

**Algorithm 13.** Batch Stochastic Gradient Descent (Algorithm 26 in Handout 3) with learning rate  $\eta_t \geq 0$  and batch size m can be used to solve (2.1) by using the update rule as

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, and set  $a \tilde{\mathcal{S}}_m = \{z_i; i \in \mathcal{J}^{(t)}\}.$
- (2) compute

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

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

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

Remark 15. The implementation of the SGD variants (Algorithms 26, 38, 43, 49 in Handout 3)) is straightforward based on Algorithm 13 and (2.4).

## 3. Full 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), the Bayesian estimator of h := h(w) can be computed as the posterior expectation of w given the data  $S_n$ 

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

Remark 17. Monte Carlo integration aims at approximating (1.3), by using Central Limit Theorem or Law of Large Numbers arguments as  $\hat{h} \approx \mathbb{E}_f(h(w)|z_{1:n})$  where

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

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

Remark 18. Stochastic gradient Langevin dynamics (SGLD) algorithm is able to approximately produce samples from the posterior distribution (1.2) of parameters w given the available data  $z_{1:n}$ . That allows to recover the whole posterior distribution (1.2) (hence account for the uncertainty in parameters) and approximate posterior expectations (1.3) by averaging out (3.1) according to the Monte Carlo integration (Remark 17).

Remark 19. Stochastic gradient Langevin dynamics (SGLD) algorithm is able to generate a sample approximately distributed according to the posterior distribution (1.2). That allows to recover the whole posterior distribution (hence account for uncertainty in parameters) and approximate posterior expectations based on the Monte Carlo integration (Remark 17).

Note 20. SGLD relies on injecting the 'right' amount of noise to a standard stochastic gradient optimization recursion (2.2), such that, as the stepsize  $\eta_t$  properly reduces, the produced chain  $\{w^{(t+1)}\}$  converges to samples that could have been drawn from the true posterior distribution.

**Algorithm 21.** Stochastic Gradient Langevin Dynamics (SGLD) with learning rate  $\eta_t > 0$ , batch size m, and temperature  $\tau > 0$  is

- (1) Generate a random set  $\mathcal{J}^{(t)} \subseteq \{1,...,n\}^m$  of m indices from 1 to n with or without replacement, and set  $a \, \tilde{\mathcal{S}}_m = \{z_i; i \in \mathcal{J}^{(t)}\}.$
- (2) Compute

(3.2) 
$$w^{(t+1)} = w^{(t)} + \eta_t \left( \frac{n}{m} \sum_{i \in J^{(t)}} \nabla \log f(z_i|w) + \nabla \log f(w) \right) + \sqrt{\eta_t} \sqrt{\tau} \epsilon_t$$

where  $\epsilon_t \stackrel{\text{iid}}{\sim} N(0,1)$ .

(3) Terminate if a termination criterion is satisfied; E.g.,  $t \leq T_{max}$  for a prespecified  $T_{max} > 0$ .

Remark 22. The first few iterations from Algorithm 21 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.

Remark 23. The output of SGLD (Algorithm 21)  $\{w^{(t)}\}$  includes the generated values of w produced during the last few iterations of the running algorithm (aka the end tail of the generated chain). Page 4 Created on 2023/02/03 at 16:36:57 by Georgios Karagiannis

Remark 24. SGLD (Algorithm 21) generates as output a random chain  $\{w^{(t)}\}$  that is approximately distributed according to a distribution with density such as

(3.3) 
$$f_{\tau}\left(w|z_{1:n}\right) \propto \exp\left(\frac{1}{\tau} \prod_{i=1}^{n} f\left(z_{i}|w\right) f\left(w\right)\right)$$

(3.4) 
$$\propto \exp\left(\frac{1}{\tau}L_n(w)f(w)\right)$$

under regularity conditions. Conditions 25 on the learning rate are rather inevitable and should be satisfied.

Condition 25. Regarding the learning rate (or gain)  $\{\eta_t\}$  should satisfy conditions

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

Remark 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.

Remark 27. SGLD for  $\tau = 1$  approximately simulates from the posterior (1.2).

Remark 28. The popular learning rates  $\{\eta_t\}$  in Remark 9 in Handout 2 satisfy Condition 25 and hence can be used in SGD too. Once parametrized,  $\eta_t$  can be tuned based on pilot runs using a reasonably small number of data.

Remark 29. (Mathematically speaking) The stochastic chain in (3.2) can be viewed as a discretization of the continuous-time Langevin diffusion described by the stochastic differential equation

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

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

is the unique invariant distribution of (3.5), and that the distributions of W(t) converge rapidly to  $f^*$  as  $t \to \infty$ .

Remark 30. (Heuristically speaking) In the initial phase of running, the stochastic gradient noise will dominate the injected noise  $\epsilon_t$  and the algorithm will imitate an efficient SGD Algorithm 10 -but this is until  $\eta_t$  or  $\nabla \log (L_n(w))$  become small enough. In the later phase of running, the injected noise  $\epsilon_t$  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

<sup>&</sup>lt;sup>1</sup>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  $0 \le s \le t$ .

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.

Remark 31. One can argue that, the output of SGLD is also an "almost" minimizer of the empirical risk for large enough t. A draw from the Gibbs distribution (3.6) is approximately a minimizer 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.

Remark 32. To guarantee the algorithm to work it is important for the step sizes  $\eta_t$  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  $\eta_t$  constant once it has decreased below a critical level.

Remark 33. 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  $\widehat{h_T(w)} \to \mathcal{E}_f(h(w)|z_{1:n})$  based on LLN arguments.

Remark 34. Another more efficient estimator for the expectation (1.3) is the weighted arithmetic average

(3.8) 
$$\widehat{h(w)} = \sum_{t=T_0+1}^{T} \frac{\eta_t}{\sum_{t=T_0+1}^{T} \eta_{t'}} h\left(w^{(t)}\right)$$

Because the step size  $\eta_t$  decreases, the mixing rate of the chain  $\{w^{(t)}\}$  decreases as well and the simple sample average (3.7) will overemphasize the tail end of the sequence where there is higher correlation among the samples resulting in higher variance in the estimator.

Remark 35. Certain dimensions may have a vastly larger curvature leading to much bigger gradients. In this case a symmetric preconditioning matrix  $P_t > 0$  can transform all dimensions to the same scale; this is similar to the SGD case in Handout 3. Hence the update (3.2) becomes

(3.9) 
$$w^{(t+1)} = w^{(t)} + \eta_t P_t \left( \frac{n}{m} \sum_{i \in J^{(t)}} \nabla \log f(z_i | w) + \nabla \log f(w) \right) + \sqrt{\eta_t} \sqrt{\tau} P_t^{\frac{1}{2}} \epsilon_t$$

where  $P_t^{\frac{1}{2}}$  is such that  $P_t^{\frac{1}{2}} \left( P_t^{\frac{1}{2}} \right)^{\top} = P_t$ .

We continue the Example 33 in Handout 1, and Example 8 in Handout 2. 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

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

Because SGD (Algorithm 13) and SGLD (Algorithm 21) 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

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

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

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

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

$$\log (f(z_{i} = (x_{i}, y_{i}) | w)) = -\frac{1}{2} \log (2\pi) - \frac{1}{2} \gamma - \frac{1}{2} (y_{j} - x_{i}^{\top} \beta)^{2} \exp(-\gamma)$$

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

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

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

<sup>&</sup>lt;sup>2</sup>Code is available from https://github.com/georgios-stats/Machine\_Learning\_and\_Neural\_Networks\_III\_ Epiphany\_2023/tree/main/Lecture\_handouts/code/04.Stochastic\_gradient\_Langevine\_dynamics

Hence for the log sampling PDF we have

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

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

To implement SGD (Algorithm 13) and SGLD (Algorithm 21), 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  $(\beta, \sigma^2)$ .

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

- 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  $\tau$  but the same batch sizes m=100 and compared it against the exact posterior densities. We observe that increasing the temperature  $\tau$  may increase the variation of the produced chain. We can use a large  $\tau$  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  $\tau=1$ .

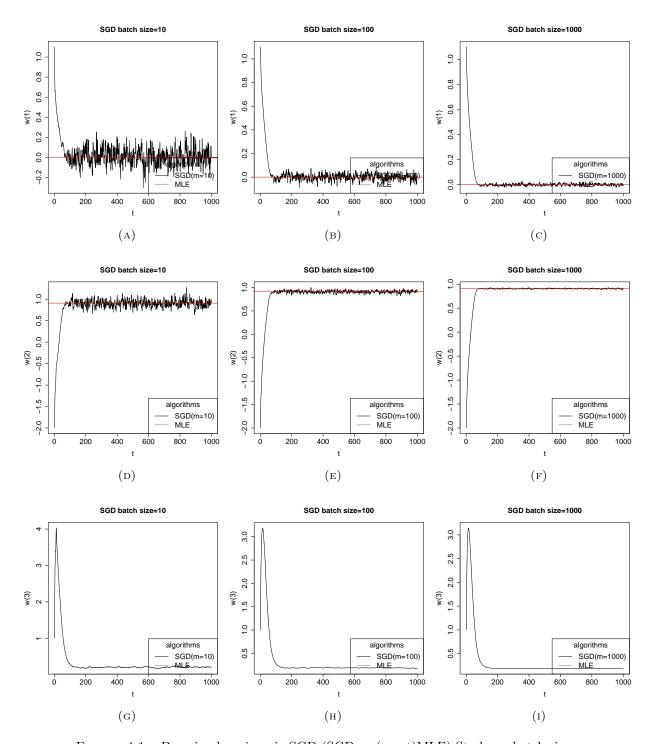


FIGURE 4.1. Bayesian learning via SGD (SGD vs (exact)MLE) Study on batch size m.

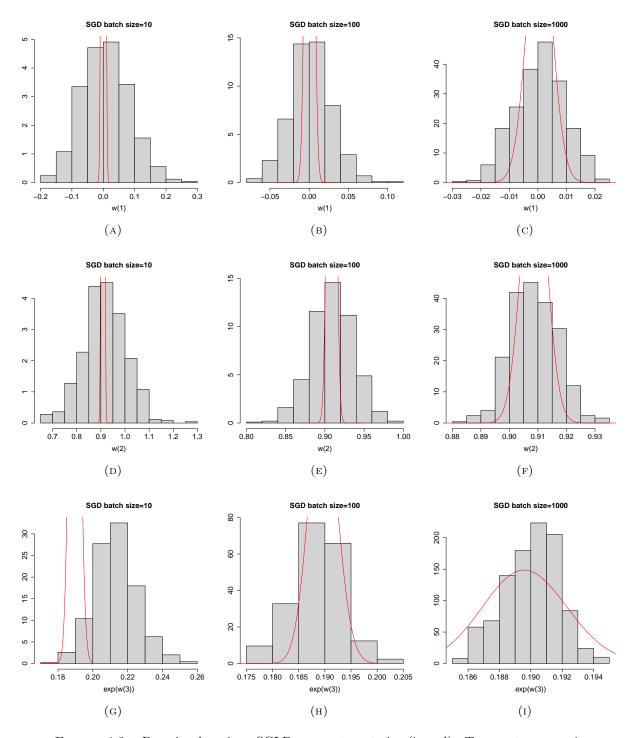


FIGURE 4.2. Bayesian learning: SGLD vs exact posterior (in red). Temperature  $\tau=1.$  Study on batch size m

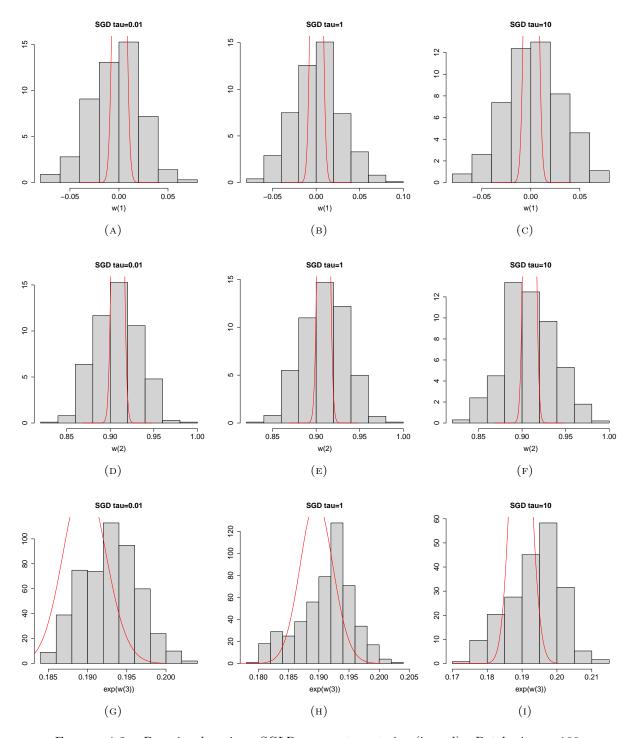


FIGURE 4.3. Bayesian learning: SGLD vs exact posterior (in red). Batch size = 100. Study on  $\tau$