

---

# Stochastic Gradient MCMC with Repulsive Forces

---

**Victor Gallego**  
ICMAT-CSIC  
victor.gallego@icmat.es

**David Rios Insua**  
ICMAT-CSIC  
david.rios@icmat.es

## Abstract

We propose a unifying view of two different families of Bayesian inference algorithms, SG-MCMC and SVGD. We show that SVGD plus a noise term can be framed as a multiple chain SG-MCMC method. Instead of treating each parallel chain independently from others, the proposed algorithm implements a repulsive force between particles, avoiding collapse. Experiments in both synthetic distributions and real datasets show the benefits of the proposed scheme.

## 1 Introduction

Bayesian computation lies at the heart of many machine learning models in both academia and industry. Thus, it is of major importance to develop more efficient approximation techniques that tackle the intractable integrals that arise in large scale Bayesian problems. Recent developments in Bayesian techniques applied to large scale datasets or deep models include variational approaches such as *Automatic Differentiation Variational Inference* (ADVI) [1] and *Stein Variational Gradient Descent* (SVGD) [2], or sampling approaches such as *Stochastic Gradient Markov Chain Monte Carlo* (SG-MCMC) [3]. While variational techniques enjoy faster computations than the latter approaches, they rely on optimizing a family of posterior approximates that may not contain the actual posterior distribution. In this work, we draw on a similitude between SG-MCMC and SVGD in order to propose an efficient sampling algorithm.

Any competing MCMC algorithm should verify the following list of properties: *scalability* (we resort to SG-MCMC methods since at each iteration they may be approximated to just require a minibatch of the dataset), *convergence to the true posterior* and *flexibility* (since we provide a parametric formulation of the transition kernel it is possible to adapt other methods such as *Hamiltonian Monte Carlo* [4] or the Nosé-Hoover thermostat method [5]).

Our contributions are summarized in what follows:

- We provide a unifying hybrid scheme of SG-MCMC and SVGD algorithms, satisfying the previous list of requirements.
- Based on it, we can develop new SG-MCMC schemes that include repulsive forces between particles.

## 2 Background and related work

The celebrated work of [3] proposed a general formulation of a continuous-time Markov process that converges to a target distribution  $\pi \propto \exp(-H(\mathbf{z}))$  with  $\mathbf{z} \in \mathbb{R}^d$ . It is based on the Euler-Maruyama discretization of the generalized Langevin dynamics:

$$\mathbf{z}_{t+1} \leftarrow \mathbf{z}_t - \epsilon_t [(\mathbf{D}(\mathbf{z}_t) + \mathbf{Q}(\mathbf{z}_t))\nabla H(\mathbf{z}_t) + \mathbf{\Gamma}(\mathbf{z}_t)] + \mathcal{N}(\mathbf{0}, 2\epsilon_t \mathbf{D}(\mathbf{z}_t)), \quad (1)$$

where  $\epsilon_t$  is the stepsize,  $\mathbf{D}(\mathbf{z})$  is a diffusion matrix,  $\mathbf{Q}(\mathbf{z})$  is a curl matrix and  $\mathbf{\Gamma}(\mathbf{z})_i = \sum_{j=1}^d \frac{\partial}{\partial \mathbf{z}_j} (\mathbf{D}_{ij}(\mathbf{z}) + \mathbf{Q}_{ij}(\mathbf{z}))$  is a correction term. Hence, to obtain any valid SG-MCMC algorithm

we simply have to choose the dimensionality of  $\mathbf{z}$  (i.e., if we augment the space with auxiliary variables as in HMC),  $\mathbf{D}$  and  $\mathbf{Q}$ . For instance, the popular Stochastic gradient Langevin dynamics (SGLD) [6] is obtained when  $\mathbf{D} = \mathbf{I}$  and  $\mathbf{Q} = \mathbf{0}$ . In addition, the Hamiltonian variant can also be recovered if we augment the state space with a  $d$ -dimensional momentum term  $\bar{\mathbf{z}} = (\mathbf{z}, \mathbf{p})$ . Then, we set  $\mathbf{D} = \mathbf{0}$  and  $\mathbf{Q} = \begin{pmatrix} \mathbf{0} & -\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{pmatrix}$ .

On the other hand, SVGD [2] frames posterior sampling as an optimization process, in which a set of  $K$  particles  $\{\mathbf{z}_i\}_{i=1}^K$  is evolved iteratively via the velocity field  $\mathbf{z}_{i,t+1} \leftarrow \mathbf{z}_{i,t} + \epsilon \phi(\mathbf{z}_{i,t})$ . Assuming  $q$  is the particle distribution at iteration  $t$  and  $q_{[\epsilon\phi]}$  is the distribution after an update ( $t+1$ ), then, the optimal choice of the velocity field  $\phi$  can be framed into the optimization problem  $\phi^* = \arg \max_{\phi \in \mathcal{F}} \{-\frac{d}{d\epsilon} \text{KL}(q_{[\epsilon\phi]} \| p)\}$ , i.e.,  $\phi$  is chosen so as to maximize the decreasing rate on the KL divergence between the particle distribution and the target. When  $\mathcal{F}$  is a RKHS, [2] showed that the optimal velocity field is given by

$$\mathbf{z}_{i,t+1} \leftarrow \mathbf{z}_{i,t} - \epsilon_t \frac{1}{K} \sum_{j=1}^K [k(\mathbf{z}_{j,t}, \mathbf{z}_{i,t}) \nabla H(\mathbf{z}_{j,t}) + \nabla_{\mathbf{z}_{j,t}} k(\mathbf{z}_{j,t}, \mathbf{z}_{i,t})], \quad (2)$$

where the RBF kernel  $k(\mathbf{z}, \mathbf{z}') = \exp(-\frac{1}{h} \|\mathbf{z} - \mathbf{z}'\|^2)$  is typically adopted. Note that the gradient term  $\nabla_{\mathbf{z}_{j,t}} k(\mathbf{z}_{j,t}, \mathbf{z}_{i,t})$  acts as a repulsive force that prevents particles from collapsing. [7] started to consider similitudes between SG-MCMC and SVGD, though in this work we propose the first hybrid scheme between both methods.

### 3 Proposed scheme

We use the framework introduced in [3] in an augmented state space  $\mathbf{z} = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K)$  to obtain a valid posterior sampler that runs multiple Markov chains with interaction. This multiple particle version of SG-MCMC is given by the following equation

$$\mathbf{z}_{t+1} \leftarrow \mathbf{z}_t - \epsilon_t [(\mathbf{K} + \mathbf{Q})\nabla + \mathbf{\Gamma}] + \boldsymbol{\eta}_t, \quad \boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, 2\epsilon_t \mathbf{K}). \quad (3)$$

Now,  $\mathbf{z}_t = (\mathbf{z}_{1,t} \dots \mathbf{z}_{K,t})^\top$  is a  $Kd$ -dimensional vector defined by the concatenation of  $K$  particles.  $\nabla \in \mathbb{R}^{K \times d \times 1}$  so that  $(\nabla)_{i,:} = \nabla H(\mathbf{z}_{i,t})^1$ .  $\mathbf{K} \in \mathbb{R}^{Kd \times Kd}$  is an expansion of the  $D$  matrix, accounting for distance between the particles and  $\mathbf{Q} \in \mathbb{R}^{Kd \times Kd}$  might be used if a Hamiltonian variant is to be adopted.  $\mathbf{\Gamma}$  is the correction term from the framework of [3]. Note that  $\mathbf{K}$ ,  $\mathbf{Q}$  and  $\mathbf{\Gamma}$  can depend on the state  $\mathbf{z}_t$  (an example will be given below), but we do not make it explicit in order to ease the notation.

Recall that, in matrix notation, the update rule for SVGD can be expressed as

$$\bar{\mathbf{z}}_{t+1} \leftarrow \bar{\mathbf{z}}_t - \frac{\epsilon_t}{K} (\bar{\mathbf{K}}\bar{\nabla} + \bar{\mathbf{\Gamma}}) \quad (4)$$

where  $\bar{\mathbf{K}} \in \mathbb{R}^{K \times K}$  so that  $(\bar{\mathbf{K}})_{ij} = k(\mathbf{z}_i, \mathbf{z}_j)$ ,  $\bar{\nabla} \in \mathbb{R}^{K \times d}$  and  $\bar{\mathbf{z}}_t \in \mathbb{R}^{K \times d}$ . Casting the later matrix as a tensor  $\nabla \in \mathbb{R}^{K \times d \times 1}$  and the former one as a tensor  $\mathbf{K} \in \mathbb{R}^{(K \times d) \times (K \times d)}$  by broadcasting along the second and fourth axes, we may associate  $\mathbf{K}$  with SG-MCMC's diffusion matrix  $\mathbf{D}$  over a  $Kd$ -dimensional space. Appendix A shows the precise definition of the matrix  $\mathbf{K}$ .

From this perspective, Eq. (4) (SVGd) can be seen as a special case of Eq. (3) when the curl matrix  $\mathbf{Q} = \mathbf{0}$  and the noise term is added. We shall refer to this perturbed variant of SVGD as *Parallel SGLD plus repulsion* (SGLD+R):

$$\mathbf{z}_{t+1} \leftarrow \mathbf{z}_t - \frac{\epsilon_t}{K} (\mathbf{K}\nabla + \mathbf{\Gamma}) + \boldsymbol{\eta}_t, \quad \boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, 2\epsilon_t \mathbf{K}/K). \quad (5)$$

Since matrix  $\mathbf{K}$  is definite positive (it was constructed from the RBF kernel), we may now use the key result from [3] (Theorem 1) to derive the following property:

**Proposition 1.** *SGLD+R (or its general form, Eq. (3)) has  $\pi(\mathbf{z}) = \prod_{k=1}^K \pi(\mathbf{z}_k)$  as stationary distribution, and the proposed discretizations are asymptotically exact as  $\epsilon_t \rightarrow 0$ .*

<sup>1</sup>Though  $\nabla \in \mathbb{R}^{Kd \times 1}$  to allow multiplication by  $\mathbf{K} + \mathbf{Q}$ , we reshape it as  $\nabla \in \mathbb{R}^{K \times d \times 1}$  to better illustrate how it is defined.

Shown that SVGD plus a noise term can be framed as a SG-MCMC method, we may now explore the design spaces of the  $\mathbf{K}$  and  $\mathbf{Q}$  matrices. However, for the rest of the paper we resort to the case  $\mathbf{Q} = \mathbf{0}$  (i.e., we will just experimentally study SGLD with repulsion between multiple particles).

Algorithm 1 shows how to set it up. Finally, our proposed method is amenable to parallelization, since the mini-batch setting from SG-MCMC can be adopted (see Appendix B).

## 4 Experiments

This Section describes the experiments developed to empirically test the proposed scheme. Code will be released at <https://github.com/vicgalle/sgmcmc-force> and we rely on Tensorflow-probability [8] as the main package.

**Synthetic distributions.** We test our proposed scheme in the following distributions.

- **Mixture of Exponentials (MoE).** Two exponential distributions with different scale parameters  $\lambda_1 = 1.5, \lambda_2 = 0.5$  and mixture proportions  $\pi_1 = 1/3, \pi_2 = 2/3$ .
- **Mixture of 2D Gaussians (MoG).** A grid of  $3 \times 3$  equally distributed isotropic 2D Gaussians, see Figure 2(d) for the density plot.

We compare two sampling methods, SGLD with  $K$  parallel chains, and our proposed scheme, SGLD+R. Note that the only difference in these two sampling algorithms is that for the former  $\mathbf{K} = \mathbf{I}$  whereas the latter accounts for repulsion between particles. Table 1 reports the effective sampling size metrics for each method using  $K = 10$  particles. Note that while ESS/s are similar, the repulsive forces in SGLD+R makes for a more efficient exploration, resulting in much lower estimation errors. Figures 1 and 2 seem to confirm this fact. In addition, even when increasing the number of particles  $K$ , SGLD+R achieves lower errors than SGLD (see Fig. 3).

| Distribution | ESS   |              | ESS/s       |             | Error of $\mathbb{E}[X]$ |             |
|--------------|-------|--------------|-------------|-------------|--------------------------|-------------|
|              | SGLD  | SGLD+R       | SGLD        | SGLD+R      | SGLD                     | SGLD+R      |
| MoE          | 44.3  | <b>59.1</b>  | 51.5        | <b>61.0</b> | 0.39                     | <b>0.14</b> |
| MoG          | 151.3 | <b>169.5</b> | <b>36.3</b> | 32.5        | 1.42                     | <b>1.19</b> |

Table 1: Results for the two synthetic distributions task

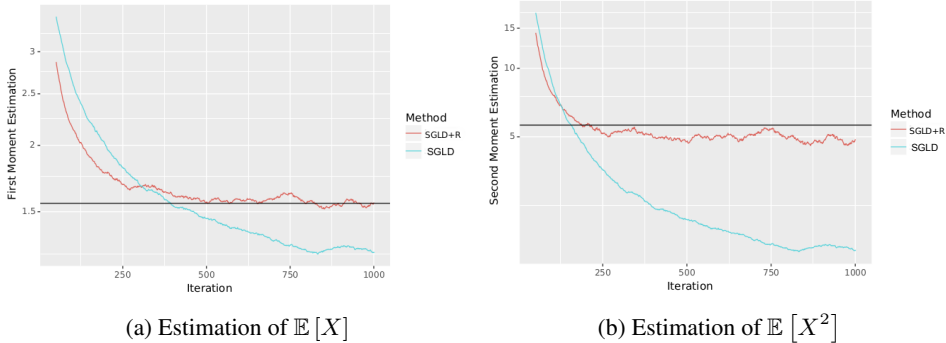


Figure 1: Evolution of estimation during the MoE experiment. 10 particles are used and black line depicts the exact value to be estimated

**Bayesian Neural Network.** We test the proposed scheme in a suite of regression tasks using a feed-forward neural network with 50 hidden units and ReLU activations. The datasets are taken from the UCI repository. We use minibatches of size 100, but while we use the same experimental setting as [2], we simply use SGD instead of Adagrad, since it is not clear that Proposition 1 can be extended to the non-SGD case. As before, we compare SGLD and SGLD+R, reporting average root mean squared error and log-likelihood over a test set in Table 2. We observe that SGLD+R typically outperforms SGLD. During the experiments, we noted that in order to reduce computation time, during the last half of training we could disable the repulsion between particles without incurring in performance cost.

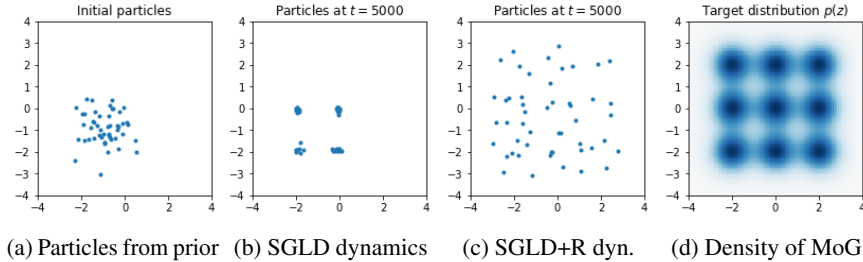


Figure 2: Evolution of the particles during the MoG experiment

| Dataset | Avg. Test RMSE    |                                     | Avg. Test LL       |                                      |
|---------|-------------------|-------------------------------------|--------------------|--------------------------------------|
|         | SGLD              | SGLD+R                              | SGLD               | SGLD+R                               |
| Boston  | $2.392 \pm 0.018$ | <b><math>2.295 \pm 0.017</math></b> | $-2.551 \pm 0.018$ | $-2.575 \pm 0.007$                   |
| Kin8nm  | $0.104 \pm 0.001$ | $0.104 \pm 0.001$                   | $0.826 \pm 0.005$  | $0.831 \pm 0.006$                    |
| Naval   | $0.008 \pm 0.000$ | $0.008 \pm 0.000$                   | $3.379 \pm 0.011$  | <b><math>3.428 \pm 0.019</math></b>  |
| Protein | $4.810 \pm 0.003$ | <b><math>4.794 \pm 0.003</math></b> | $-2.991 \pm 0.000$ | <b><math>-2.987 \pm 0.001</math></b> |
| Wine    | $0.522 \pm 0.004$ | <b><math>0.514 \pm 0.004</math></b> | $-0.765 \pm 0.008$ | <b><math>-0.750 \pm 0.007</math></b> |
| Yacht   | $0.942 \pm 0.015$ | <b><math>0.894 \pm 0.029</math></b> | $-1.211 \pm 0.020$ | $-1.172 \pm 0.026$                   |

Table 2: Results for the BNN experiments

## 5 Conclusions and further work

This paper shows how to generate new SG-MCMC methods consisting in multiple chains plus repulsion between the particles. Instead of the naïve parallelization, in which a particle from a chain is agnostic to the others, we showed how it is possible to adapt another method from the literature, SVGD, in order to account for better exploration of the space, avoiding collapse between particles. Further work shall explore different matrices  $\mathbf{K}$  and  $\mathbf{Q}$  in order to further accelerate the sampling process.

**Acknowledgments** Authors thank Roi Naveiro from the ICMAT Datalab group for interesting discussions. Authors also thank the anonymous reviewers for their suggestions. VG acknowledges support from grant FPU16-05034. DRI is grateful to the MINECO MTM2014-56949-C3-1-R project and the AXA-ICMAT Chair in Adversarial Risk Analysis. All authors acknowledge support from the Severo Ochoa Excellence Programme SEV-2015-0554.

## References

- [1] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- [2] Qiang Liu and Dilin Wang. Stein variational gradient descent: A general purpose bayesian inference algorithm. In *Advances In Neural Information Processing Systems*, pages 2378–2386, 2016.
- [3] Yi-An Ma, Tianqi Chen, and Emily Fox. A complete recipe for stochastic gradient mcmc. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2917–2925. Curran Associates, Inc., 2015.
- [4] Tianqi Chen, Emily Fox, and Carlos Guestrin. Stochastic gradient hamiltonian monte carlo. In *International Conference on Machine Learning*, pages 1683–1691, 2014.
- [5] Nan Ding, Youhan Fang, Ryan Babbush, Changyou Chen, Robert D Skeel, and Hartmut Neven. Bayesian sampling using stochastic gradient thermostats. In *Advances in Neural Information Processing Systems*, pages 3203–3211, 2014.
- [6] Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 681–688, 2011.
- [7] Qiang Liu. Stein variational gradient descent as gradient flow. In *Advances in Neural Information Processing Systems*, pages 3115–3123, 2017.



$$\nabla_z \log \pi(z|\mathcal{D}) \approx \nabla_z \log p(z) + \frac{N}{|\Omega|} \sum_{i \in \Omega} \nabla_z \log p(D_i|z).$$

## C Experiment details

### Synthetic distributions

For the Mixture of Exponentials experiment, the pdf used is

$$p(z) = \sum_{i=1}^2 \pi_i \lambda_i \exp(-\lambda_i z).$$

Note that the exact value of the first and second moments can be computed using

$$\mathbb{E}[z^n] = \sum_{i=1}^2 \pi_i \frac{n!}{\lambda_i^n}$$

with  $n \in \mathbb{N}$ . Since  $z > 0$ , in order to use the proposed scheme, we reparameterize using the log function. The pdf of the transformation  $y = \log(z)$  can be computed using

$$p(y) = p(\log^{-1}(y)) |D \log^{-1}(y)|.$$

For the Mixture of Gaussians experiment, we set  $\Sigma = \text{diag}(0.1, 0.1)$  and place each of the nine Gaussians centered at each each of the following points  $\{(-2, -2), (-2, 0), (-2, 2), (0, -2), (0, 0), (0, 2), (2, -2), (2, 0), (2, 2)\}$ .

For the computation of the error of  $\mathbb{E}[X]$  in Table 1, we sample for 500 iterations after discarding the first 500 iterations as burn-in, and we collected samples every 10 iterations to reduce correlation between samples. For the MoE case we used 10 particles whereas for the MoG task we used 20 particles due to the increased number of modes.

### Bayesian Neural Network

Learning rate was chosen from a grid  $\{1e-5, \dots, 1e-3\}$  validated on another fold. The number of iterations was set to 2000 in every experiment. As before, to make predictions we collect samples every 10 iterations after a burn-in period. 20 particles were used for each of the tested datasets.

## D Additional results

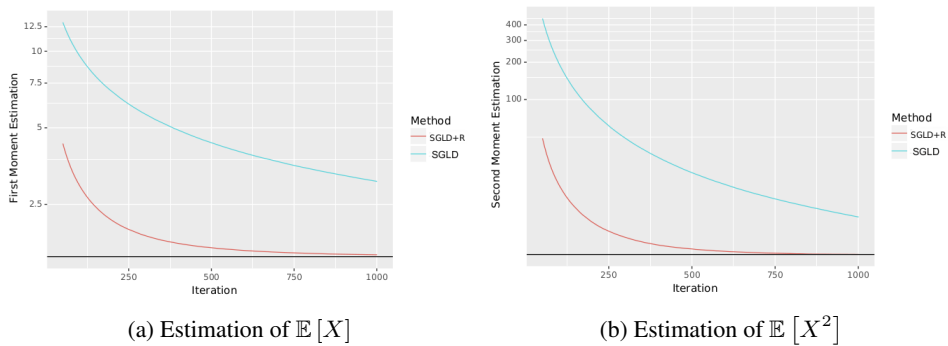


Figure 3: Evolution of estimation during the MoE experiment. 100 particles are used and black line depicts the exact value to be estimated