

Intro and Goals

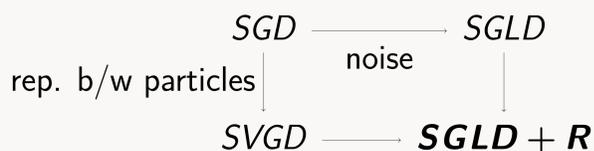
- Recent developments in Bayesian techniques applied to large scale datasets or deep models include **variational approaches** such as *Automatic Differentiation Variational Inference* (ADVI) [1] and *Stein Variational Gradient Descent* (SVGD) [2], or **sampling approaches** such as *Stochastic Gradient Markov Chain Monte Carlo* (SG-MCMC) [3].
- Can we bridge the gap between variational and sampling methods?
Yes, we propose an hybrid between SGLD and SVGD!

Background

- SG-MCMC [3]
 - Choose state space $\mathbf{z} \in \mathbb{R}^d$ and target distribution $\pi \propto \exp(-H(\mathbf{z}))$.
 - Choose suitable diffusion $\mathbf{D}(\mathbf{z})$ and curl $\mathbf{Q}(\mathbf{z})$ matrices.
 - Discretize the generalized Langevin dynamics:
$$\mathbf{z}_{t+1} \leftarrow \mathbf{z}_t - \epsilon_t [(\mathbf{D}(\mathbf{z}_t) + \mathbf{Q}(\mathbf{z}_t))\nabla H(\mathbf{z}_t) + \mathbf{\Gamma}(\mathbf{z}_t)] + \boldsymbol{\eta}_t,$$
where $\boldsymbol{\eta}_t$ is some carefully chosen Gaussian noise.
 - Stochastic Gradient Langevin Dynamics (SGLD):
 $\mathbf{D} = \mathbf{I}$ and $\mathbf{Q} = \mathbf{0}$.
 - Hamiltonian variant (HMC): $\bar{\mathbf{z}} = (\mathbf{z}, \mathbf{p})$.
 $\mathbf{D} = \mathbf{0}$ and $\mathbf{Q} = \begin{pmatrix} \mathbf{0} & -\mathbf{I} \\ \mathbf{I} & \mathbf{0} \end{pmatrix}$
- 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
$$\mathbf{z}_{i,t+1} \leftarrow \mathbf{z}_{i,t} - \epsilon_t \frac{1}{K} \sum_{j=1}^K [k(\mathbf{z}_{j,t}, \mathbf{z}_{i,t})\nabla H(\mathbf{z}_{j,t}) + \nabla_{\mathbf{z}_{j,t}} k(\mathbf{z}_{j,t}, \mathbf{z}_{i,t})],$$
where the RBF kernel $k(\mathbf{z}, \mathbf{z}') = \exp(-\frac{1}{h}\|\mathbf{z} - \mathbf{z}'\|^2)$ is typically adopted. This velocity field is chosen so as to to maximize the decreasing rate on the KL divergence between the particle distribution and the target.

Proposed scheme

Instead of using K parallel chains without interactions, we propose



Parallel SGLD plus repulsion (SGLD+R):

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

where \mathbf{K} is a permuted block-diagonal matrix such that for each block $(\mathbf{K})_{i,j} = k(\mathbf{z}_{j,t}, \mathbf{z}_{i,t})$. (i.e., instead of identity or diagonal diffusion matrix as in SGLD, we use a block-diagonal matrix accounting for distances between particles)

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

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

Experiments

Synthetic distributions:

- Mixture of Exponentials (MoE).
- Mixture of 2D Gaussians (MoG).

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

Table: Results for the two synthetic distributions task

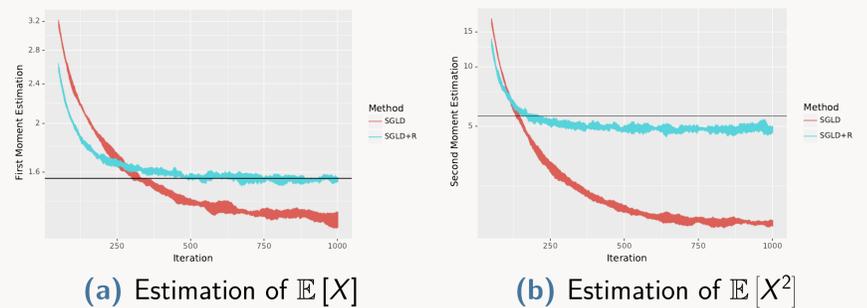


Figure: Evolution of estimation during the MoE experiments (5 simulations). 10 particles are used for each sim. and black line depicts the exact value to be estimated

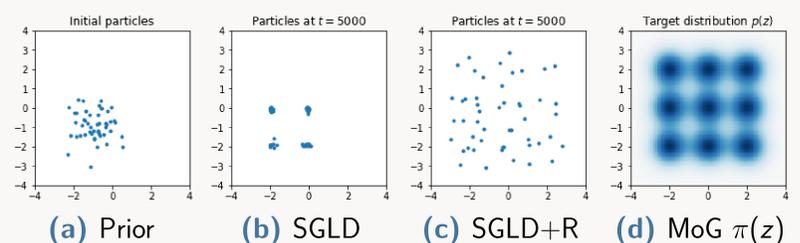


Figure: Evolution of the particles during the MoG experiment

Bayesian Neural Network:

Feed-forward neural network over some regression tasks from the UCI datasets.

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

Conclusions and Further Work

- We showed how to generate new SG-MCMC methods consisting in multiple chains plus repulsion between the particles.
- Repulsion between particles improves exploration of the space, avoiding particle collapse. Plus, we may collect much more samples than with SVGD.
- Explore different matrices \mathbf{K} and \mathbf{Q} in order to further accelerate the sampling process

$$\mathbf{z}_{t+1} \leftarrow \mathbf{z}_t - \epsilon_t [(\mathbf{K} + \mathbf{Q})\nabla + \mathbf{\Gamma}] + \boldsymbol{\eta}_t. \quad (1)$$

References

- 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.
- Qiang Liu and Dilin Wang. Stein Variational Gradient Descent: A General Purpose Bayesian Inference Algorithm. *In Advances In Neural Information Processing Systems*, 2016.
- Yi-An Ma, Tianqi Chen, and Emily Fox. A Complete Recipe for Stochastic Gradient MCMC. *In Advances in Neural Information Processing Systems*. 2015.